

# Estimation And Selection Via Absolute Penalized Convex Minimization And Its Multistage Adaptive Applications

Jian Huang and Cun-Hui Zhang<sup>1</sup>

University of Iowa and Rutgers University

*Abstract:* The  $\ell_1$ -penalized method, or the Lasso, has emerged as an important tool for the analysis of large data sets. Many important results have been obtained for the Lasso in linear regression which have led to a deeper understanding of high-dimensional statistical problems. In this article, we consider a class of weighted  $\ell_1$ -penalized estimators for convex loss functions of a general form, including the generalized linear models. We study the estimation, prediction, selection and sparsity properties of the weighted  $\ell_1$ -penalized estimator in sparse, high-dimensional settings where the number of predictors  $p$  can be much larger than the sample size  $n$ . Adaptive Lasso is considered as a special case. A multistage method is developed to apply an adaptive Lasso recursively. We provide  $\ell_q$  oracle inequalities, a general selection consistency theorem, and an upper bound on the dimension of the Lasso estimator. Important models including the linear regression, logistic regression and log-linear models are used throughout to illustrate the applications of the general results.

*Running title:* Absolute penalized convex minimization

*Key words:* Variable selection, penalized estimation, oracle inequality, generalized linear models, selection consistency, sparsity.

## 1 Introduction

High-dimensional data arise in many diverse fields of scientific research. For example, in genetic and genomic studies, more and more large data sets are being generated with rapid advances in biotechnology, where the total number of variables  $p$  is larger than the sample size  $n$ . Fortunately, statistical analysis is still possible for a

---

<sup>1</sup>Jian Huang's research is partially supported by NIH Grants CA120988, CA142774 and NSF Grant DMS 0805670. Cun-Hui Zhang's research is partially supported by NSF Grants DMS 0604571, DMS 0804626 and NSA Grant H98230-09-1-0006.

substantial subset of such problems with a sparse underlying model where the number of important variables is much smaller than the sample size. A fundamental problem in the analysis of such data is to find reasonably accurate sparse solutions that are easy to interpret and can be used for the prediction and estimation of covariable effects. The  $\ell_1$ -penalized method, or the Lasso [Tib96, CDS98], has emerged as an important approach to finding such solutions in sparse, high-dimensional statistical problems.

In the last few years, considerable progress has been made in understanding the theoretical properties of the Lasso in  $p \gg n$  settings. Most results have been obtained for linear regression models with a quadratic loss. [GR04] studied the prediction performance of the Lasso in high-dimensional least squares regression. [MB06] showed that, for neighborhood selection in the Gaussian graphical models, under a neighborhood stability condition on the design matrix and certain additional regularity conditions, the Lasso is selection consistent even when  $p \rightarrow \infty$  at a rate faster than  $n$ . [ZY06] formalized the neighborhood stability condition in the context of linear regression as a strong irrepresentable condition. [CT07] derived an upper bound for the  $\ell_2$  loss for the estimation of regression coefficients with a closely related Dantzig selector under a condition on the number of nonzero coefficients and a uniform uncertainty principle on the design matrix. Similar results have been obtained for the Lasso. For example, upper bounds for the  $\ell_q$  loss of the Lasso estimator has been established by [BTW07] for  $q = 1$ , [ZH08] for  $q \in [1; 2]$ , [MY09] for  $q = 2$ , [BRT09] for  $q \in [1; 2]$ , and [Zha09, YZ10] for general  $q \geq 1$ . For convex minimization methods beyond linear regression, [vdG08] studied the Lasso in high-dimensional generalized linear models (GLM) and obtained prediction and  $\ell_1$  estimation error bounds. [NRWY10] studied penalized M-estimators with a general class of regularizers, including an  $\ell_2$  error bound for the Lasso in GLM under a restricted convexity and other regularity conditions.

Theoretical studies of the Lasso have revealed that it may not perform well for the purpose of variable selection, since its required irrepresentable condition is not properly scaled in the number of relevant variables. In a number of simulation studies, the Lasso has shown weakness in variable selection when the number of nonzero regression coefficients increases. As a remedy, a number of proposals have been introduced in the literature, including concave penalized LSE [FL01, Zha10a],

adaptive Lasso [Zou06], and stepwise regression [Zha11]. Although extensions of the concave penalized LSE is beyond the scope of this paper, adaptive Lasso is studied here as a weighted Lasso with estimated weights. When the number of predictors  $p$  is fixed, [Zou06] proved that the adaptive Lasso has the asymptotic oracle property. In linear regression models. [HMZ08] showed that the oracle property continues to hold for the adaptive Lasso in  $p \gg n$  settings under an adaptive irrerepresentable and other regularity conditions. [MB07] suggested using the Lasso as the initial estimator for the adaptive Lasso or even a multi-step adaptive Lasso. The one-step method of [ZL08], designed to approximate penalized estimators with concave penalties, can be also viewed as adaptive Lasso.

In this article, we consider a class of weighted  $\ell_1$ -penalized estimators with a convex loss function. This class includes the Lasso, adaptive Lasso and multistage recursive application of an adaptive Lasso in generalized linear models as special cases. We study the estimation, prediction, selection and sparsity properties of the weighted  $\ell_1$ -penalized estimator based on a convex loss in sparse, high-dimensional settings where the number of predictors  $p$  can be much larger than the sample size  $n$ . The main contributions of this work are as follows.

- We extend the existing theory for the unweighted Lasso from linear regression to more general convex loss function.
- We develop a multistage method with recursive applications of an adaptive Lasso and provide sharper risk bound than those for unweighted Lasso.
- We apply our general results to a number of important special cases, including the linear, logistic and log-linear regression models.

This article is organized as follows. In Section 2 we describe a general formulation of the absolute penalized minimization problem with a convex loss, along with two basic inequalities and a number of examples. In Section 3 we develop oracle inequalities for the weighted Lasso estimator for general quasi star-shaped loss functions and an  $\ell_2$  bound on the prediction error. In Section 4 we develop sharper oracle inequalities for multistage recursive applications of an adaptive Lasso. In Section 5 we derive sufficient conditions for selection consistency. In Section 6 we provide an upper bound on the dimension of the Lasso estimator. Concluding remarks are given in Section 7. All proofs are provided in an appendix.

## 2 Absolute penalized convex minimization

### 2.1 Definition and the KKT conditions

We consider a general convex loss function of the form

$$\ell(\beta) = \psi(\beta) - \langle \beta, z \rangle, \quad (1)$$

where  $\psi(\beta)$  is a known convex function,  $z$  is observed and  $\beta$  is unknown. Unless otherwise stated, the inner product space is  $\mathbb{R}^p$ , so that  $\{z, \beta\} \subset \mathbb{R}^p$  and  $\langle \beta, z \rangle = \beta'z$ . Our analysis of (1) requires certain smoothness of the function  $\psi(\beta)$  in terms of its differentiability. In what follows, such smoothness assumptions are always explicitly described by invoking the derivative of  $\psi$ . For any  $v = (v_1, \dots, v_p)'$ , we use  $\|v\|$  to denote a general norm of  $v$  and  $|v|_q$  the  $\ell_q$  norm  $(\sum_j |v_j|^q)^{1/q}$ , with  $|v|_\infty = \max_j |v_j|$ . Let  $\widehat{w} \in \mathbb{R}^p$  be a (possibly estimated) weight vector with nonnegative elements  $\widehat{w}_j, 1 \leq j \leq p$ , and  $\widehat{W} = \text{diag}(\widehat{w})$ . The weighted absolute penalized estimator, or weighted Lasso, is defined as

$$\widehat{\beta} = \arg \min_{\beta} \left\{ \ell(\beta) + \lambda |\widehat{W}\beta|_1 \right\}. \quad (2)$$

Here we focus on the case where  $\widehat{W}$  is diagonal. In linear regression, [TT11] considered non-diagonal, predetermined  $\widehat{W}$  and derived an algorithm for computing the solution paths.

A vector  $\widehat{\beta}$  is a global minimizer in (2) if and only if the negative gradient at  $\widehat{\beta}$  satisfies the Karush-Kuhn-Tucker (KKT) conditions,

$$g = -\dot{\ell}(\widehat{\beta}) = z - \dot{\psi}(\widehat{\beta}), \quad \begin{cases} g_j = \widehat{w}_j \lambda \text{sgn}(\widehat{\beta}_j) & \text{if } \widehat{\beta}_j \neq 0 \\ g_j \in \widehat{w}_j \lambda [-1, 1] & \text{all } j, \end{cases} \quad (3)$$

where  $\dot{\ell}(\beta) = (\partial/\partial\beta)\ell(\beta)$  and  $\dot{\psi}(\beta) = (\partial/\partial\beta)\psi(\beta)$ . Since the KKT conditions are necessary and sufficient for (2), results on the performance of  $\widehat{\beta}$  can be viewed as analytical consequences of (3).

The estimator (2) includes the  $\ell_1$ -penalized estimator, or the Lasso, with the choice  $\widehat{w}_j = 1, 1 \leq j \leq p$ . A careful study of the (unweighted) Lasso in general convex minimization (1) is by itself an interesting and important problem. Our work includes the Lasso as a special case since  $\widehat{w}_j = 1$  is allowed in all our theorems.

In practice, unequal  $\widehat{w}_j$  arise in many ways. In adaptive Lasso [Zou06], a decreasing function of a certain initial estimator of  $\beta_j$  is used as the weight  $\widehat{w}_j$  to remove the bias of the Lasso. In [FL01, ZL08, Zha10b], the weights  $\widehat{w}_j$  are computed iteratively with  $\widehat{w}_j = \dot{\rho}_\lambda(\widehat{\beta}_j)$ , where  $\dot{\rho}_\lambda(t) = (d/dt)\rho_\lambda(t)$  with a suitable concave penalty function  $\rho_\lambda(t)$ . This is also designed to remove the bias of the Lasso, since the concavity of  $\rho_\lambda(t)$  guarantees smaller weight for larger  $\widehat{\beta}_j$ . In Section 4, we provide results on the improvements of this weighted Lasso over the standard Lasso. In linear regression, [Zha10b] gave suitable conditions under which this iterative algorithm provides smaller weights  $\widehat{w}_j$  for most large  $\beta_j$ . Such nearly unbiased methods are expected to produce better results than the Lasso when a significant fraction of nonzero  $|\beta_j|$  are of the order  $\lambda$  or larger. Regardless of the computational methods, the results in this paper demonstrate the benefits of using data dependent weights in a general class of problems with convex losses.

Unequal weights may also arise for computational reasons. The Lasso with  $\widehat{w}_j = 1$  is expected to perform similarly to weighted Lasso with data dependent  $1 \leq \widehat{w}_j \leq C_0$ , with a fixed  $C_0$ . However, the weighted Lasso is easier to compute since  $\widehat{w}_j$  can be determined as a part of an iterative algorithm. For example, in a gradient descent algorithm, one may take larger steps and stop the computation as soon as the KKT conditions (3) are attained for any weights satisfying  $1 \leq \widehat{w}_j \leq C_0$ .

The weight function  $\widehat{w}_j$  can be also used to standardize the penalty level, for example with  $\widehat{w}_j = \{\ddot{\psi}_{jj}(\widehat{\beta})\}^{1/2}$ , where  $\ddot{\psi}_{jj}(\beta)$  is the  $j$ -th diagonal element of the Hessian matrix of  $\psi(\beta)$ . When  $\psi(\beta)$  is quadratic, for example in linear regression,  $\widehat{w}_j$  does not depend on  $\widehat{\beta}$ . However, in other convex minimization problems, such weights need to be computed iteratively.

Finally, in certain applications, the effects of a certain set  $S_*$  of variables are of primary interest, so that penalization of  $\widehat{\beta}_{S_*}$ , and thus the resulting bias, should be avoided. This leads to “semi-penalized” estimators with  $\widehat{w}_j = 0$  for  $j \in S_*$ , for example, with  $\widehat{w}_i = I\{j \notin S_*\}$ .

## 2.2 Basic inequalities, prediction, and Bregman divergence

Let  $\beta^*$  denote a target vector for  $\beta$ . In high-dimensional models, the performance of an estimator  $\widehat{\beta}$  is typically measured by its proximity to a target under conditions

on the sparsity of  $\beta^*$  and the size of the negative gradient  $-\dot{\ell}(\beta^*) = z - \dot{\psi}(\beta^*)$ . For  $\ell_1$ -penalized estimators, such results are often derived from the KKT conditions (3) via certain basic inequalities, which are direct consequences of the KKT conditions and have appeared in different forms in the literature, for example, in the papers cited in the Introduction. Let  $D(\beta, \beta^*) = \ell(\beta) - \ell(\beta^*) - \langle \dot{\ell}(\beta^*), \beta - \beta^* \rangle$  be the Bregman divergence [Bre67] and consider its symmetrized version [NN07]

$$\Delta(\beta, \beta^*) = D(\beta, \beta^*) + D(\beta^*, \beta) = \langle \beta - \beta^*, \dot{\psi}(\beta) - \dot{\psi}(\beta^*) \rangle. \quad (4)$$

Since  $\psi$  is convex,  $\Delta(\beta, \beta^*) \geq 0$ . Two basic inequalities below provide upper bounds for the symmetrized Bregman divergence  $\Delta(\widehat{\beta}, \beta^*)$ . The sparsity of  $\beta^*$  is measured by a weighted  $\ell_1$  norm of  $\beta^*$  in the first one and by the number of zero entries in the second one.

Let  $S$  be any set of indices satisfying  $S \supseteq \{j : \beta_j^* \neq 0\}$  and let  $S^c$  be the complement of  $S$  in  $\{1, \dots, p\}$ . We shall refer to  $S$  as the sparse set. Let  $W = \text{diag}(w)$  for a possibly unknown vector  $w \in \mathbb{R}^p$  with elements  $w_j \geq 0$ . Define

$$z_0^* = |\{z - \dot{\psi}(\beta^*)\}_S|_\infty, \quad z_1^* = |W_{S^c}^{-1}\{z - \dot{\psi}(\beta^*)\}_{S^c}|_\infty, \quad (5)$$

$$\Omega_0 = \{\widehat{w}_j \leq w_j \forall j \in S\} \cap \{w_j \leq \widehat{w}_j \forall j \in S^c\}, \quad (6)$$

where for any  $p$ -vector  $v$  and set  $A$ ,  $v_A = (v_j : j \in A)'$ . Here and in the sequel  $M_{AB}$  denotes the  $A \times B$  subblock of a matrix  $M$  and  $M_A = M_{AA}$ .

**Lemma 1** (i) Let  $\beta^*$  be a target vector. In the event  $\Omega_0 \cap \{|(z - \dot{\psi}(\beta^*))_j| \leq \widehat{w}_j \lambda \forall j\}$ ,

$$\Delta(\widehat{\beta}, \beta^*) \leq 2\lambda |\widehat{W}\beta^*|_1 \leq 2\lambda |W\beta^*|_1. \quad (7)$$

(ii) For any target vector  $\beta^*$  and  $S \supseteq \{j : \beta_j^* \neq 0\}$ , the error  $h = \widehat{\beta} - \beta^*$  satisfies

$$\begin{aligned} \Delta(\beta^* + h, \beta^*) + (\lambda - z_1^*) |W_{S^c} h_{S^c}|_1 &\leq \langle h_S, g_S - \{z - \dot{\psi}(\beta^*)\}_S \rangle \\ &\leq (|w_S|_\infty \lambda + z_0^*) |h_S|_1 \end{aligned} \quad (8)$$

in  $\Omega_0$  for a certain negative gradient vector  $g$  satisfying  $|g_j| \leq \widehat{w}_j \lambda$ . Consequently, in  $\Omega_0 \cap \{(|w_S|_\infty \lambda + z_0^*)/(\lambda - z_1^*) \leq \xi\}$ ,  $h \neq 0$  belongs to the sign-restricted cone  $\mathcal{C}_-(\xi, S) = \{b \in \mathcal{C}(\xi, S) : b_j(\dot{\psi}(\beta + b) - \dot{\psi}(\beta))_j \leq 0 \forall j \in S^c\}$ , where

$$\mathcal{C}(\xi, S) = \{b \in \mathbb{R}^p : |W_{S^c} b_{S^c}|_1 \leq \xi |b_S|_1 \neq 0\}. \quad (9)$$

**Remark 2.1** *Sufficient conditions are given in Subsection 3.2 for  $\{|(z - \dot{\psi}(\beta^*))_j| \leq \hat{w}_j \lambda \forall j\}$  to hold with high probability in generalized linear models. See Lemma 2, Remarks 3.3 and 3.4 and Examples 3.2, 3.3, and 3.4.*

A useful feature of Lemma 1 is the explicit statements of the monotonicity of the basic inequality in the weights. By Lemma 1 (ii), it suffices to study the analytical properties of the penalized criterion with the error  $h = \hat{\beta} - \beta^*$  in the sign-restricted cone, provided that the event  $(|w_S|_\infty \lambda + z_0^*)/(\lambda - z_1^*) \leq \xi$  has large probability. However, unless  $\mathcal{C}_-(\xi, S)$  is specified, we will consider the larger cone in (9) in order to simplify the analysis. The choices of the target vector  $\beta^*$ , the sparse set  $S \supseteq \{j : \beta_j^* \neq 0\}$ , weight vector  $\hat{w}$  and its bound  $w$  are quite flexible. The main requirement is that  $\{|S|, z_0^*, z_1^*\}$  should be small. In linear regression or generalized linear models, we may conveniently consider  $\beta^*$  as the vector of true regression coefficients under a probability measure  $P_{\beta^*}$ . However,  $\beta^*$  can also be a sparse version of a true  $\beta$ , e.g.  $\beta_j^* = \beta_j I\{|\beta_j| \geq \tau\}$  for a threshold value  $\tau$  under  $P_\beta$ .

The upper bound in Lemma 1 (i) gives the so called “slow rate” of convergence for the Bregman divergence. In Section 3, we provide “fast rate” of convergence for the Bregman divergence via oracle inequalities for  $|h_S|_1$  in (8). The symmetrized Bregman divergence  $\Delta(\hat{\beta}, \beta^*)$  has the interpretations as the regret in prediction error in linear regression, the symmetrized Kullback-Leibler (KL) divergence in generalized linear models (GLM) and density estimation, and a spectrum loss for the graphical Lasso, as shown in examples below.

**Example 2.1 (*Linear regression*)** *Consider the linear regression model*

$$y_i = \sum_{j=1}^p x_{ij} \beta_j + \varepsilon_i, \quad i = 1, \dots, n, \quad (10)$$

where  $y_i$  is the response variable,  $x_{ij}$  are predictors or design variables, and  $\varepsilon_i$  is the error term. Let  $y = (y_1, \dots, y_n)'$  and let  $X$  be the design matrix whose  $i$ th row is  $x^i = (x_{i1}, \dots, x_{ip})$ . The estimator (2) is a weighted Lasso with  $\psi(\beta) = |X\beta|_2^2/(2n)$  and  $z = X'y/n$  in (1). For predicting a vector  $\tilde{y}$  with  $E_{\beta^*}[\tilde{y}|X, y] = X\beta^*$ ,

$$\begin{aligned} n\Delta(\hat{\beta}, \beta^*) &= |X\hat{\beta} - X\beta^*|_2^2 \\ &= E_{\beta^*}[|\tilde{y} - X\hat{\beta}|_2^2|X, y] - \min_{\delta(X, y)} E_{\beta^*}[|\tilde{y} - \delta(X, y)|_2^2|X, y] \end{aligned}$$

is the regret of using the linear predictor  $X\hat{\beta}$  compared with the optimal predictor. See [GR04] for several implications of (7).

**Example 2.2 (Logistic regression)** We observe  $(X, y) \in \mathbb{R}^{n \times (p+1)}$  with independent rows  $(x^i, y_i)$ , where  $y_i \in \{0, 1\}$  are binary response variables with

$$P_\beta(y_i = 1|x^i) = \pi_i(\beta) = \exp(x^i\beta)/(1 + \exp(x^i\beta)), \quad 1 \leq i \leq n. \quad (11)$$

The loss function (1) is the average negative log-likelihood

$$\ell(\beta) = \psi(\beta) - z'\beta \quad \text{with} \quad \psi(\beta) = \sum_{i=1}^n \frac{\log(1 + \exp(x^i\beta))}{n}, \quad z = X'y/n. \quad (12)$$

Thus, (2) is a weighted  $\ell_1$  penalized MLE. For probabilities  $\{\pi', \pi''\} \subset (0, 1)$ , the KL information is  $K(\pi', \pi'') = \pi' \log(\pi'/\pi'') + (1 - \pi') \log\{(1 - \pi')/(1 - \pi'')\}$ . Since  $\dot{\psi}(\beta) = \sum_{i=1}^n x^i \pi_i(\beta)/n$  and  $\text{logit}(\pi_i(\beta^*)) - \text{logit}(\pi_i(\beta)) = x^i(\beta^* - \beta)$ , (4) gives

$$\Delta(\beta, \beta^*) = \frac{1}{n} \sum_{i=1}^n \left\{ K(\pi_i(\beta^*), \pi_i(\beta)) + K(\pi_i(\beta), \pi_i(\beta^*)) \right\}.$$

Thus,  $\Delta(\beta^*, \beta)$  is the symmetrised KL-divergence.

**Example 2.3 (GLM).** The GLM contains the linear and logistic regression models as special cases. We observe  $(X, y) \in \mathbb{R}^{n \times (p+1)}$  with rows  $(x^i, y_i)$ . Suppose that conditionally on  $X$ ,  $y_i$  are independent under  $P_\beta$  with

$$y_i \sim f(y_i|\theta_i) = \exp\left(\frac{\theta_i y_i - \psi_0(\theta_i)}{\sigma^2} + \frac{c(y_i, \sigma)}{\sigma^2}\right), \quad \theta_i = x^i \beta. \quad (13)$$

Let  $f_{(n)}(y|X, \beta) = \prod_{i=1}^n f(y_i|x^i\beta)$ . The loss function can be written as a normalized negative likelihood  $\ell(\beta) = (\sigma^2/n) \log f_{(n)}(y|X, \beta)$  with  $z = X'y/n$  and  $\psi(\beta) = \sum_{i=1}^n \{\psi_0(x^i\beta) + c(y_i, \sigma)\}/n$ . The KL divergence is

$$D(f_n(\cdot|X, \beta^*) \| f_n(\cdot|X, \beta)) = E_{\beta^*} \log \left( \frac{f_{(n)}(y|X, \beta^*)}{f_{(n)}(y|X, \beta)} \right).$$

The symmetrized Bregman divergence can be written as

$$\Delta(\hat{\beta}, \beta^*) = \frac{\sigma^2}{n} \left\{ D(f_{(n)}(\cdot|X, \beta^*) \| f_{(n)}(\cdot|X, \hat{\beta})) + D(f_{(n)}(\cdot|X, \hat{\beta}) \| f_{(n)}(\cdot|X, \beta^*)) \right\}. \quad (14)$$



**Example 2.4 (Nonparametric density estimation)** Although the focus of this paper is on regression models, here we illustrate that  $\Delta(\widehat{\beta}, \beta^*)$  is the symmetrised KL divergence in the context of nonparametric density estimation. Suppose the observations  $y = (y_1, \dots, y_n)'$  are iid from  $f(\cdot|\beta) = \exp\{\langle \beta, T(\cdot) \rangle - \psi(\beta)\}$  under  $P_\beta$ , where  $T(\cdot) = (u_j(\cdot), j \leq p)'$  with certain basis functions  $u_j(\cdot)$ . Let the loss function  $\ell(\beta)$  in (1) be the average negative log-likelihood  $n^{-1} \sum_{i=1}^n \log f(y_i|\beta)$  with  $z = n^{-1} \sum_{i=1}^n T(y_i)$ . Since  $E_\beta T(y_i) = \dot{\psi}(\beta)$ , the KL divergence is

$$D(f(\cdot|\beta^*)||f(\cdot|\beta)) = E_{\beta^*} \log \left( \frac{f(y_i|\beta^*)}{f(y_i|\beta)} \right) = \psi(\beta) - \psi(\beta^*) - \langle \beta - \beta^*, \dot{\psi}(\beta^*) \rangle.$$

Again, the symmetrised KL divergence between the target density  $f(\cdot|\beta^*)$  and the estimated density  $f(\cdot|\widehat{\beta})$  is

$$\Delta(\beta, \beta^*) = D(f(\cdot|\beta^*)||f(\cdot|\widehat{\beta})) + D(f(\cdot|\widehat{\beta})||f(\cdot|\beta^*)). \quad (15)$$

[vdG08] pointed out that for this example, the natural choices of the basis functions  $u_j$  and weights  $w_j$  satisfy  $\int u_j d\nu = 0$  and  $w_k^2 = \int u_k^2 d\nu$ .

**Example 2.5 (Graphical Lasso)** Suppose we observe  $X \in \mathbb{R}^{n \times p}$  and would like to estimate the precision matrix  $\beta = (EX'X/n)^{-1} \in \mathbb{R}^{p \times p}$ . In the graphical Lasso, (1) is the length normalized negative likelihood with  $\psi(\beta) = -\log \det \beta$ ,  $z = -X'X/n$ , and  $\langle \beta, z \rangle = -\text{trace}(\beta z)$ . Since  $\dot{\psi}(\beta) = E_\beta z = -\beta^{-1}$ , we find

$$\Delta(\beta, \beta^*) = \text{trace}((\widehat{\beta} - \beta^*)((\beta^*)^{-1} - \widehat{\beta}^{-1})) = \sum_{j=1}^p (\lambda_j - 1)^2 / \lambda_j, \quad (16)$$

where  $(\lambda_1, \dots, \lambda_p)$  are the eigenvalues of  $(\beta^*)^{-1/2} \widehat{\beta} (\beta^*)^{-1/2}$ . In graphical Lasso, the diagonal elements are typically not penalized. Consider  $\widehat{w}_{jk} = I\{j \neq k\}$ , so that the penalty for the off-diagonal elements are uniformly weighted. Since Lemma 1 requires  $|(z - \dot{\psi}(\beta^*))_{jk}| \leq \widehat{w}_{jk} \lambda$ ,  $\beta^*$  is taken to match  $X'X/n$  on the diagonal and the true  $\beta^o$  in correlations. Let  $S = \{(j, k) : \beta_{jk}^o \neq 0, j \neq k\}$ . In the event  $\max_{j \neq k} |z_{jk} - \beta_{jk}^*| \leq \lambda$ , Lemma 1 (i) gives  $\|(\beta^*)^{-1/2} \widehat{\beta} (\beta^*)^{-1/2} - I_{p \times p}\|_2 = o(1)$  under the condition  $|S| \lambda \max_{j \neq k} |\beta_{jk}^*| = o(1)$ , where  $\|\cdot\|_2$  is the spectrum norm. [RBLZ08] proved the consistency of the graphical Lasso under similar conditions with a different analysis.

### 3 Oracle inequalities

In this section, we extract upper bounds for the estimation error  $\widehat{\beta} - \beta^*$  from the basic inequality (8). Since (8) is monotone in the weights, the oracle inequalities are sharper when the weights  $\widehat{w}_j$  are smaller in  $S = \{j : \beta_j^* \neq 0\}$  and larger in  $S^c$ .

We say that a function  $\phi(b)$  defined in  $\mathbb{R}^p$  is quasi star-shaped if  $\phi(tb)$  is continuous and non-decreasing in  $t \in [0, \infty)$  for all  $b \in \mathbb{R}^p$  and  $\lim_{b \rightarrow 0} \phi(b) = 0$ . All seminorms are quasi star-shaped. The sublevel sets  $\{b : \phi(b) \leq t\}$  of a quasi star-shaped function are all star-shaped. For  $0 \leq \eta^* \leq 1$  and any pair of quasi star-shaped functions  $\phi_0(b)$  and  $\phi(b)$ , define

$$F(\xi, S; \phi_0, \phi) = \inf \left\{ \frac{\Delta(\beta^* + b, \beta^*) e^{\phi_0(b)}}{|b_S|_1 \phi(b)} : b \in \mathcal{C}(\xi, S), \phi_0(b) \leq \eta^* \right\}, \quad (17)$$

where  $\Delta(\beta, \beta^*)$  is as in (4). We refer to  $F(\xi, S; \phi_0, \phi)$  as a general invertibility factor (GIF) over the cone (9). The GIF plays a crucial role in developing the error bounds for  $\widehat{\beta} - \beta^*$ . It extends the squared compatibility constant [vdGB09] and the weak and sign-restricted cone invertibility factors [YZ10] from the linear regression model with  $\phi_0(\cdot) = 0$  to more general model (1) and from  $\ell_q$  norms to general  $\phi(\cdot)$ . They are all closely related to the restricted eigenvalues [BRT09, Kol09] as we will discuss in Subsection 3.1.

The basic inequality (8) implies that the symmetrized Bregman divergence  $\Delta(\widehat{\beta}, \beta^*)$  is no greater than a linear function of  $|h_S|_1$ , where  $h = \widehat{\beta} - \beta^*$ . If  $\Delta(\widehat{\beta}, \beta^*)$  is no smaller than a linear function of the product  $|h_S|_1 \phi(h)$ , then an upper bound for  $\phi(h)$  exists. Since the symmetrized Bregman divergence (4) is approximately quadratic,  $\Delta(\widehat{\beta}, \beta^*) \approx h' \ddot{\psi}(\beta^*) h$ , in a neighborhood of  $\beta^*$ , this is reasonable when  $h = \widehat{\beta} - \beta^*$  is not too large and  $\ddot{\psi}(\beta^*)$  is invertible in the cone. A suitable factor  $e^{\phi_0(b)}$  in (17) forces the computation of this lower bound in a proper neighborhood of  $\beta^*$ .

We first provide a set of general oracle inequalities.

**Theorem 1** *Let  $\{z_0^*, z_1^*\}$  be as in (5) with  $S \supseteq \{j : \beta_j^* \neq 0\}$ ,  $\Omega_0$  in (6),  $0 \leq \eta \leq \eta^* \leq 1$ , and  $\{\phi_0(b), \phi(b)\}$  be a pair of quasi star-shaped functions. Let  $\phi_{1,S}(b) = |b_S|_1 / |S|$ . In the event*

$$\Omega_1 = \Omega_0 \cap \left\{ \frac{|w_S|_\infty \lambda + z_0^*}{\lambda - z_1^*} \leq \xi, \frac{|w_S|_\infty \lambda + z_0^*}{F(\xi, S; \phi_0, \phi)} \leq \eta e^{-\eta} \right\}, \quad (18)$$

the following oracle inequalities hold:

$$\phi_0(\widehat{\beta} - \beta^*) \leq \eta, \quad \phi(\widehat{\beta} - \beta^*) \leq \frac{e^\eta(|w_S|_\infty \lambda + z_0^*)}{F(\xi, S; \phi_0, \phi)}, \quad (19)$$

$$\Delta(\widehat{\beta}, \beta^*) + (\lambda - z_1^*)|W_{S^c}(\widehat{\beta} - \beta^*)_{S^c}|_1 \leq \frac{e^\eta(|w_S|_\infty \lambda + z_0^*)^2 |S|}{F(\xi, S; \phi_0, \phi_{1,S})}. \quad (20)$$

**Remark 3.1** *Sufficient conditions are given in Subsection 3.2 for (18) to hold with high probability. See Lemma 2, Remarks 3.3 and 3.4 and Examples 3.2, 3.3, and 3.4.*

The oracle inequalities in Theorem 1 control both the estimation error in terms of  $\phi_0(\widehat{\beta} - \beta^*)$  and the prediction error in terms of the symmetrized Bregman divergence  $\Delta(\widehat{\beta}, \beta^*)$  discussed in Section 2. Since they are based on (17) in the intersection of the cone and the unit ball  $\{b : \phi_0(b) \leq 1/e\}$ , they are different from typical results in a small-ball analysis based on the Taylor expansion of  $\psi(\beta)$  at  $\beta = \beta^*$ . Theorem 1 does allow  $\phi_0(\cdot) = 0$  with  $F(\xi, S; \phi_0, \phi_0) = \infty$  and  $\eta = 0$  in linear regression.

### 3.1 The Hessian and related quantities

We describe the relationship between the GIF (17) and the Hessian of the convex function  $\psi(\cdot)$  in (1) and examine cases where the quasi star-shaped functions  $\phi_0(\cdot)$  and  $\phi(\cdot)$  are familiar seminorms. Throughout, we assume that  $\psi(\beta)$  is twice differentiable. Let  $\ddot{\psi}(\beta)$  be the Hessian of  $\psi(\beta)$  and  $\Sigma^* = \ddot{\psi}(\beta^*)$ .

The GIF (17) can be simplified if for a certain nonnegative-definite matrix  $\Sigma$ ,

$$\Delta(\beta^* + b, \beta^*) e^{\phi_0(b)} \geq \langle b, \Sigma b \rangle, \quad \forall b \in \mathcal{C}(\xi, S), \quad \phi_0(b) \leq \eta^*. \quad (21)$$

Since  $\Delta(\beta^* + h, \beta^*) = \int_0^1 \langle h, \ddot{\psi}(\beta^* + th)h \rangle dt$  by (4), (21) is a smoothness condition on the Hessian when  $\Sigma = \Sigma^*$ . In what follows,  $\Sigma = \Sigma^*$  is allowed in all statements unless otherwise stated. Under (21), (17) is bounded from below by the simple GIF,

$$F_0(\xi, S; \phi) = \inf_{b \in \mathcal{C}(\xi, S)} \frac{\langle b, \Sigma b \rangle}{|b_S|_1 \phi(b)}. \quad (22)$$

In linear regression,  $F_0(\xi, S; \phi)$  is the square of the compatibility factor for  $\phi(b) = \phi_{1,S}(b) = |b_S|_1/|S|$  [vdG07] and the cone invertibility factor for  $\phi(b) = \phi_q(b) = |b|_q/|S|^{1/q}$  [YZ10]. They are both closely related to the restricted isometry property

(RIP) [CT05], the sparse Rieze condition (SRC) [ZH08], and the restricted eigenvalue [BRT09]. Extensive discussion of these quantities can be found in [BRT09, vdGB09, YZ10]. The following corollary is an extension of an oracle inequality of [YZ10] for the linear regression model.

**Corollary 1** *Let  $\eta \leq \eta^* \leq 1$ . Suppose (21) holds. Then, in the event*

$$\Omega_0 \cap \{|w_S|_\infty \lambda + z_0^* \leq \min(\xi(\lambda - z_1^*), \eta e^{-\eta} F_0(\xi, S; \phi_0))\},$$

(19) and (20) hold with  $F(\xi, S; \phi_0, \phi)$  replaced by the simpler  $F_0(\xi, S; \phi)$  in (22). In particular, in the same event,

$$\phi_0(h) \leq \eta, \quad |h|_q \leq \frac{e^\eta (|w_S|_\infty \lambda + z_0^*) |S|^{1/q}}{F_0(\xi, S; \phi_q)}, \quad \forall q > 0, \quad (23)$$

with  $\phi_q(b) = |b|_q / |S|^{1/q}$  and  $h = \hat{\beta} - \beta^*$ , and with  $\phi_{1,S}(b) = |b_S|_1 / |S|$ ,

$$e^{-\eta} h' \Sigma h \leq \Delta(\hat{\beta}, \beta^*) \leq \frac{e^\eta (|w_S|_\infty \lambda + z_0^*)^2 |S|}{F_0(\xi, S; \phi_{1,S})} - (\lambda - z_1^*) |W_{S^c} h_{S^c}|_1. \quad (24)$$

Here the only differences between the general model (1) and linear regression ( $\phi_0(b) = 0$ ) are the extra factor  $e^\eta$  with  $\eta \leq 1$ , the extra constraint  $|w_S|_\infty \lambda + z_0^* \leq \eta e^{-\eta} F_0(\xi, S; \phi_0)$ , and the extra condition (21). Moreover, (22) explicitly expresses all conditions on  $F_0(\xi, S; \phi)$  as properties of a fixed  $\Sigma$ .

**Example 3.1 (Linear regression: oracle inequalities).** For  $\psi(\beta) = |Xb|_2^2 / (2n)$  and  $\Sigma = X'X/n$ ,  $F_0(\xi, S; \phi_q)$  is the weak cone invertibility factor [YZ10] and  $F_0^{1/2}(\xi, S; \phi_{1,S})$  is the compatibility constant [vdG07]

$$\kappa_*(\xi, S) = \inf_{b \in \mathcal{C}(\xi, S)} \frac{|S|^{1/2} |Xb|_2}{|b_S|_1 n^{1/2}} = \inf_{b \in \mathcal{C}(\xi, S)} \left( \frac{b' \Sigma b}{|b_S|_1^2 / |S|} \right)^{1/2}. \quad (25)$$

They are all closely related to the  $\ell_2$  restricted eigenvalues

$$RE_2(\xi, S) = \inf_{b \in \mathcal{C}(\xi, S)} \frac{|Xb|_2}{|b|_2 n^{1/2}} = \inf_{b \in \mathcal{C}(\xi, S)} \left( \frac{b' \Sigma b}{|b|_2^2} \right)^{1/2} \quad (26)$$

[BRT09, Kol09]. Since  $|b_S|_1^2 \leq |b|_2^2 |S|$ ,  $\kappa_*(\xi, S) \geq RE_2(\xi, S)$  [vdGB09]. For the Lasso with  $\hat{w}_j = 1$ ,

$$|\hat{\beta} - \beta^*|_2 \leq \frac{|S|^{1/2} (\lambda + z_0^*)}{SCIF_2(\xi, S)} \leq \frac{|S|^{1/2} (\lambda + z_0^*)}{F_0(\xi, S; \phi_2)} \leq \frac{|S|^{1/2} (\lambda + z_0^*)}{\kappa_*(\xi, S) RE_2(\xi, S)} \quad (27)$$

in the event  $\lambda + z_0^* \leq \xi(\lambda - z_1^*)$  [YZ10], where

$$SCIF_q(\xi, S) = \inf_{b \in \mathcal{C}_-(\xi, S)} |\Sigma b|_\infty / \phi_q(b), \quad \phi_q = |b|_q / |S|^{1/q}.$$

Thus, cone and general invertibility factors yield sharper  $\ell_2$  oracle inequalities.

The factors in the oracle inequalities in (27) do not have the same order for large  $|S|$  and certain design matrices  $X$ . Although the oracle inequality based on  $SCIF_2(\xi, S)$  is the sharpest in (27), it seems not to lead to a simple extension to the general convex minimization with (1). Thus, we settle with extensions of the second sharpest oracle inequality in (27) with  $F_0(\xi, S; \cdot)$ .

### 3.2 Oracle inequalities for the Lasso in GLM

An important special case of the general formulation is the  $\ell_1$ -penalized estimator in a generalized linear model (GLM) [MN89]. This is Example 2.3 in Subsection 2.2, where we set up the notation in (13) and gave the KL divergence interpretation to (4). The  $\ell_1$  penalized, normalized negative likelihood is

$$\ell(\beta) = \psi(\beta) - z'\beta, \quad \text{with } \psi(\beta) = C_n(y, \sigma) + \sum_{i=1}^n \frac{\psi_0(x^i \beta)}{n} \text{ and } z = \frac{X'y}{n}. \quad (28)$$

Assume that  $\psi_0$  is twice differentiable. Denote the first and second derivatives of  $\psi_0$  by  $\dot{\psi}_0$  and  $\ddot{\psi}_0$ , respectively. The gradient and Hessian are

$$\dot{\psi}(\beta) = X'\dot{\psi}_0(\theta)/n \quad \text{and} \quad \ddot{\psi}(\beta) = X'\text{diag}(\ddot{\psi}_0(\theta))X/n, \quad (29)$$

where  $\theta = X\beta$  and  $\dot{\psi}_0$  and  $\ddot{\psi}_0$  are applied to the individual components of  $\theta$ .

A crucial condition in our analysis of the Lasso in GLM is

$$\max_{i \leq n} \left| \log \left( \frac{\ddot{\psi}_0(x^i \beta^* + t)}{\ddot{\psi}_0(x^i \beta^*)} \right) \right| \leq M_1 |t|, \quad \forall M_1 |t| \leq \eta^* \quad (30)$$

where  $M_1$  and  $\eta^*$  are constants determined by  $\psi_0$ . This condition gives

$$\Delta(\beta^* + b, \beta^*) = \int_0^1 \langle b, \ddot{\psi}(\beta^* + tb)b \rangle dt \geq \int_0^1 \sum_{tM_1 |x^i b| \leq \eta^*} \frac{\ddot{\psi}_0(x^i \beta^*)(x^i b)^2}{ne^{tM_1 |x^i b|}} dt,$$

which implies the following lower bound for the GIF in (17):

$$\begin{aligned} F(\xi, S; \phi_0, \phi) &\geq \inf_{b \in \mathcal{C}(\xi, S), \phi_0(b) \leq \eta^*} \sum_{i=1}^n \frac{\ddot{\psi}_0(x^i \beta^*)(x^i b)^2}{n|b_S|_1 \phi(b)} \int_0^1 I\{tM_1|x^i b| \leq \phi_0(b)\} dt \\ &= \inf_{b \in \mathcal{C}(\xi, S), \phi_0(b) \leq \eta^*} \sum_{i=1}^n \frac{\ddot{\psi}_0(x^i \beta^*)}{n|b_S|_1 \phi(b)} \min\left(\frac{|x^i b| \phi_0(b)}{M_1}, (x^i b)^2\right), \end{aligned} \quad (31)$$

due to  $(x^i b)^2 \int_0^1 I\{tM_1|x^i b| \leq \phi_0(b)\} dt = \min\{|x^i b| \phi_0(b)/M_1, (x^i b)^2\}$ . For seminorms  $\phi_0$  and  $\phi$ , the infimum above can be taken over a fixed value of  $\phi_0(b)$  due to scale invariance. Thus, for  $\phi_0(b) = M_2|b|_2$  and seminorms  $\phi$ , the lower bound in (31) is

$$F^*(\xi, S; \phi) = \inf_{b \in \mathcal{C}(\xi, S), |b|_2=1} \sum_{i=1}^n \frac{\ddot{\psi}_0(x^i \beta^*)}{n|b_S|_1 \phi(b)} \min\left(\frac{|x^i b| M_2}{M_1}, (x^i b)^2\right). \quad (32)$$

If (30) holds with  $\eta^* = \infty$ , the convexity of  $e^{-t}$  yields (21) with

$$\phi_0(b) = \frac{M_1 \sum_{i=1}^n \ddot{\psi}_0(x^i \beta^*) |x^i b|^3}{\sum_{i=1}^n \ddot{\psi}_0(x^i \beta^*) (x^i b)^2} \leq M_1 |Xb|_\infty, \quad (33)$$

with an application of the Jensen inequality. This gives a special  $F_0(\xi, S; \phi_0)$  as

$$F_*(\xi, S) = \inf_{b \in \mathcal{C}(\xi, S)} \frac{n \langle b, \Sigma^* b \rangle^2 / (M_1 |b_S|_1)}{\sum_{i=1}^n \ddot{\psi}_0(x^i \beta^*) |x^i b|^3}. \quad (34)$$

We note that since  $|Xb|_\infty \leq |X_S|_\infty |b_S|_1 + |X_{S^c} W_{S^c}^{-1}|_\infty |W_{S^c} b_{S^c}| \leq \{|X_S|_\infty + \xi |X_{S^c} W_{S^c}^{-1}|_\infty\} |b_S|_1$  in the cone  $\mathcal{C}(\xi, S)$  in (9), for  $\phi_0(b) = M_3 |b_S|_1$  with  $M_3 = M_1 \{|X_S|_\infty + \xi |X_{S^c} W_{S^c}^{-1}|_\infty\}$ , (21) automatically implies the stronger

$$e^{-\phi_0(b)} \langle b, \Sigma^* b \rangle \leq \Delta(\beta^* + b, \beta^*) \leq e^{\phi_0(b)} \langle b, \Sigma^* b \rangle, \quad \forall b \in \mathcal{C}(\xi, S), \quad \phi_0(b) \leq \eta^*. \quad (35)$$

Under condition (30), we may also use the following large deviation inequalities to find explicit penalty levels to guarantee (18).

**Lemma 2** (i) Suppose (13) and (30) hold with certain  $\{M_1, \eta^*\}$  and the  $w_j$  in (6) are deterministic. Let  $x_j$  be the columns of  $X$ ,  $\Sigma_{ij}^*$  be the elements of  $\Sigma^* = \ddot{\psi}(\beta^*)$ . For positive constants  $\{\lambda_0, \lambda_1\}$  define  $t_j = \lambda_0 I\{j \in S\} + w_j \lambda_1 I\{j \notin S\}$ . Suppose

$$M_1 \max_{j \leq p} (|x_j|_\infty |t_j / \Sigma_{jj}^*|) \leq \eta_0 e^{\eta_0} \quad \text{and} \quad \sum_{j=1}^p \exp\left\{-\frac{nt_j^2 e^{-\eta_0}}{2\sigma^2 \Sigma_{jj}^*}\right\} \leq \frac{\epsilon_0}{2} \quad (36)$$

for certain constants  $\eta_0 \leq \eta^*$  and  $\epsilon_0 > 0$ . Then,  $P_{\beta^*} \left\{ z_0^* \leq \lambda_0, z_1^* \leq \lambda_1 \right\} \geq 1 - \epsilon_0$ .

(ii) If  $c_0 = \max_t \ddot{\psi}(t)$ , then part (i) is still valid if (30) and (36) are replaced by

$$\sum_{j=1}^p \exp \left\{ - \frac{n^2 t_j^2}{2\sigma^2 c_0 |x_j|_2^2} \right\} \leq \frac{\epsilon_0}{2}. \quad (37)$$

In particular, if  $|x_j|_2^2 = n, 1 \leq j \leq p, w_j = 1, j \notin S$  and  $\lambda_0 = \lambda_1 = \lambda$  (so  $t_j = \lambda$ ), then part (i) still holds if  $\lambda \geq \sigma \sqrt{(2c_0/n) \log(2p/\epsilon_0)}$ .

The following theorem is a consequence of Theorem 1, Corollary 1 and Lemma 2.

**Theorem 2** (i) Let  $\widehat{\beta}$  be the Lasso (2) with the loss function in (28). Let  $\beta^*$  be a target vector and  $h = \widehat{\beta} - \beta^*$ . Suppose (13) and (30) hold with certain  $\{M_1, \eta^*\}$ . Let  $F^*(\xi, S; \phi)$  be as in (32) with  $S = \{j : \beta_j^* \neq 0\}$  and a constant  $M_2$ . Let  $\eta \leq 1 \wedge \eta^*$  and  $\{\lambda, \lambda_0, \lambda_1\}$  satisfy

$$|w_S|_\infty \lambda + \lambda_0 \leq \min \left\{ \xi(\lambda - \lambda_1), \eta e^{-\eta} F^*(\xi, S; M_2 | \cdot |_2) \right\}. \quad (38)$$

Then, in the event  $\Omega_0 \cap \left\{ \max_{k=0,1} (z_k^*/\lambda_k) \leq 1 \right\}$  with the  $z_k^*$  in (5) and  $\Omega_0$  in (6),

$$\Delta(\beta^* + h, \beta^*) \leq \frac{e^\eta (|w_S|_\infty \lambda + \lambda_0)^2 |S|}{F^*(\xi, S; \phi_{1,S})}, \quad \phi(h) \leq \frac{e^\eta (|w_S|_\infty \lambda + \lambda_0)}{F^*(\xi, S; \phi)} \quad (39)$$

for all seminorms  $\phi$ . Moreover, if either (36) or (37) holds for the  $\{\lambda_0, \lambda_1\}$  and  $W$  is deterministic, then

$$P_{\beta^*} \left\{ (39) \text{ holds for all seminorms } \phi \right\} \geq P_{\beta^*}(\Omega_0) - \epsilon_0.$$

(ii) If  $\eta^* = \infty$  and (38) holds with  $F^*(\xi, S; M_2 | \cdot |_2)$  replaced by the  $F_*(\xi, S)$  in (34), then the conclusions of part (i) hold with  $F^*(\xi, S; \cdot)$  replaced by the  $F_0(\xi, S; \cdot)$  in (22). Moreover, (39) can be strengthened with the lower bound  $\Delta(\beta^* + h, \beta^*) \geq e^{-\eta} \langle h, \Sigma^* h \rangle$ .  
(iii) For any  $\eta^* > 0$ , the conclusions of part (ii) hold if  $F_*(\xi, S)$  is replaced by  $\kappa_*^2(\xi, S)/(M_3 |S|)$  in (38) with the  $M_3$  in (35).

**Remark 3.2** Since  $\phi = \phi_0$  is allowed in (39), (39) implies  $\phi_0(h) \leq \eta$  with  $\phi_0(h) = M_2 |h|_2$  in part (i) and the  $\phi_0$  in (33) in part (ii). Similarly, under the conditions of Theorem 2 (iii),  $M_3 |h_S|_1 \leq \eta \leq \eta^*$ , so that (35) holds with  $b = h = \widehat{\beta} - \beta^*$ .

**Remark 3.3** If either (36) or (37) holds for  $\{\lambda_0, \lambda_1\}$  and  $W$  is deterministic, then (39) implies  $P_{\beta^*}\{(18) \text{ holds}\} \geq P_{\beta^*}(\Omega_0) - \epsilon_0$ .

**Remark 3.4** Suppose  $\{\min_{j \notin S} w_j, \min_j \Sigma_{jj}^*\}$  are bounded away from zero,  $\{\max_{j \in S} w_j, \max_j \Sigma_{jj}^*, M_1\}$  are bounded, and  $\{1 + F_*^2(\xi, S)\}(\log p)/n \rightarrow 0$ . Then, (36) holds with  $\lambda_0 = \lambda_1 = a\sigma\sqrt{(2/n)\log(p/\epsilon_0)}$  for certain  $a \leq (1 + o(1))\max_j(\Sigma^*)_{jj}^{1/2}/w_j$ , due to  $\max\{\lambda_0, \eta, \eta_0\} \rightarrow 0+$ . Again, the conditions and conclusions of Theorem 2 “converge” to those for the linear regression as if the Gram matrix is  $\Sigma^*$ .

**Remark 3.5** In Theorem 2, the key condition (38) is weaker in parts (i) and (ii) than part (iii), although part (ii) requires  $\eta^* = \infty$ . For  $\Sigma = \Sigma^*$  and  $M_1 = M_2 \leq M_3/(1+\xi)$ ,

$$\kappa_*^2(\xi, S)/(M_3|S|) \leq \min \left\{ F_*(\xi, S), F^*(\xi, S; M_2|\cdot|_2) \right\},$$

since  $n^{-1} \sum_{i=1}^n \ddot{\psi}_0(x^i \beta^*) |x^i b|^3 / \langle b, \Sigma^* b \rangle \leq |Xb|_\infty \leq |b_S|_1 M_3/M_1$  as in the derivation of (35) and  $|b|_2 \leq (1+\xi)|b_S|_1$  in the cone (9). For the more familiar  $\kappa_*^2(\xi, S)/(M_3|S|)$ , (38) essentially requires a small  $|S|\sqrt{(\log p)/n}$ . The sharper Theorem 2 (i) and (ii) provides conditions to relax the requirement to a small  $|S|(\log p)/n$ .

**Remark 3.6** For  $\hat{w}_j = 1$ , [NRWY10] considered  $M$ -estimators under a restricted strong convexity condition. For the GLM, they considered iid sub-Gaussian  $x^i$  and used empirical process theory to bound  $\Delta(\beta^* + b, \beta^*)/\{|b|_2(|b|_2 - c_0|b|_1)\}$  from below over the cone (9) with a small  $c_0$ . Their result extends the  $\ell_2$  error bound  $|S|^{1/2}(\lambda + z_0^*)/RE_2^2(\xi, S)$  of [BRT09], while Theorem 2 extends the sharper (27) with the factor  $F_0(\xi, S; \phi_2)$ . Theorem 2 applies to both deterministic and random designs. Similar to [NRWY10], for iid sub-Gaussian  $x^i$ , empirical process theory can be used to verify (38) with  $F^*(\xi, S; M_2|\cdot|_2) \gtrsim |S|^{-1/2}$ , provided that  $|S|(\log p)/n$  is small.

**Example 3.2 (Linear regression: oracle inequalities, continuation)** For the linear regression model (10) with quadratic loss,  $\psi_0(\theta) = \theta^2/2$ , so that (30) holds with  $M_1 = 0$  and  $\eta^* = \infty$ . It follows that  $F^*(\xi, S; M_2|\cdot|_2) = \infty$  and (38) has the interpretation with  $\eta = 0+$  and  $\eta e^{-\eta} F^*(\xi, S; M_2|\cdot|_2) = \infty$ . Moreover, since  $M_1 = 0$ ,  $\eta_0 = 0+$  in (36). Thus, the conditions and conclusions of Theorem 2 “converge” to the case of linear regression as  $M_1 \rightarrow 0+$ . Suppose  $\varepsilon_i \sim N(0, \sigma^2)$  as in (13). For  $\hat{w}_j = w_j = 1$  and  $\Sigma_{jj}^* = \sum_{i=1}^n x_{ij}^2/n = 1$ , (36) holds with  $\lambda_0 = \lambda_1 = \sigma\sqrt{(2/n)\log(p/\epsilon_0)}$



and (38) holds with  $\lambda = \lambda_0(1 + \xi)/(1 - \xi)$ . The value of  $\sigma$  can be estimated iteratively using the mean residual squares [SBvdG10, SZ11]. Alternatively, cross-validation can be used to pick  $\lambda$ . For  $\phi(b) = \phi_2(b) = |b|_2/|S|^{1/2}$ , (39) matches (27) with the factor  $F_0(\xi, S; \phi_2)$ .

**Example 3.3 (Logistic regression: oracle inequalities)** The model and loss function are given in (11) and (12) respectively. Here we verify the conditions of Theorem 2. Condition (30) holds with  $M_1 = 1$  and  $\eta^* = \infty$ ; Since  $\psi_0(t) = \log(1 + e^t)$ ,

$$\frac{\ddot{\psi}_0(\theta + t)}{\ddot{\psi}_0(\theta)} = \frac{e^t(1 + e^\theta)^2}{(1 + e^{\theta+t})^2} \geq \begin{cases} e^{-|t|} & t < 0 \\ e^{-t}(1 + e^\theta)^2/(e^{-t} + e^\theta)^2 \geq e^{-|t|} & t > 0. \end{cases}$$

Since  $\max_t \ddot{\psi}(t) = c_0 = 1/4$  we can apply (37). In particular, if  $\hat{w}_j = w_j = 1 = |x_j|_2^2/n$ ,  $\lambda = \{(\xi + 1)/(\xi - 1)\} \sqrt{(\log(p/\epsilon_0))/(2n)}$  and  $\lambda\{2\xi/(\xi + 1)\}/F_*(\xi, S) \leq \eta e^{-\eta}$ , then (39) holds with at least probability  $1 - \epsilon_0$  under  $P_{\beta^*}$ . For such  $\hat{W}$  and  $X$ , an adaptive choice of the penalty level is  $\lambda = \hat{\sigma} \sqrt{(2/n) \log p}$  with  $\hat{\sigma}^2 = \sum_{i=1}^n \pi_i(\hat{\beta})\{1 - \pi_i(\hat{\beta})\}/n$ , where  $\pi_i(\beta)$  is as in Example 2.2.

**Example 3.4 (Log-linear models: oracle inequalities)** Consider counting data with  $y_i \in \{0, 1, 2, \dots\}$ . In log-linear models, it is assume that

$$E_\beta(y_i) = e^{\theta_i}, \quad \theta_i = x^i \beta, \quad 1 \leq i \leq n. \quad (40)$$

The average negative Poisson log-likelihood function is

$$\ell(\beta) = \psi(\beta) - z' \beta, \quad \psi(\beta) = \sum_{i=1}^n \frac{\exp(x^i \beta) - \log(y_i!)}{n}, \quad z = X' y / n. \quad (41)$$

Again this is a GLM. In this model,  $\psi_0(t) = e^t$ , so that (30) holds with  $M_1 = 1$  and  $\eta^* = \infty$ . Although (37) is not useful with  $c_0 = \infty$ , (36) can be used in Theorem 2.

## 4 Adaptive and multistage methods

We consider in this section an adaptive Lasso and its repeated applications, with weights recursively generated based a concave penalty function. This approach appears to provide the most appealing choice of weights both from heuristic and

theoretical standpoints. The analysis here is based on the results in Section 3 and the main idea in [Zha10b].

Let  $\rho_\lambda(t)$  be a penalty function with  $\dot{\rho}_\lambda(0+) = \lambda$ , where  $\dot{\rho}_\lambda(t) = (\partial/\partial t)\rho_\lambda(t)$ . Define

$$\kappa = \sup_{0 < t_1 < t_2} \frac{|\dot{\rho}_\lambda(t_2) - \dot{\rho}_\lambda(t_1)|}{t_2 - t_1}. \quad (42)$$

Let  $\Sigma$  be as in (21) and  $\mathcal{C}(\xi, S)$  be the cone in (9). Define

$$F_2(\xi, S) = \inf \left\{ \frac{b' \Sigma b}{|b_S|_2 |b|_2} : 0 \neq b \in \mathcal{C}(\xi, S) \right\}. \quad (43)$$

The quantity  $F_2(\xi, S)$  is slightly larger than the square of the restricted eigenvalue (26) for a design matrix  $X$  when  $\Sigma = X'X/n$ . Given  $0 < \epsilon_0 < 1$ , the components of the error vector  $z - \dot{\psi}(\beta^*)$  are sub-Gaussian if for all  $0 \leq t \leq \sigma \sqrt{(2/n) \log(4p/\epsilon_0)}$ ,

$$\mathbb{P}_{\beta^*} \left\{ |(z - \dot{\psi}(\beta^*))_j| \geq t \right\} \leq 2e^{-nt^2/(2\sigma^2)}. \quad (44)$$

This condition holds for all GLM when the components of  $X\beta^*$  are uniformly in the interior of the natural parameter space for the exponential family.

**Theorem 3** *Suppose (21) holds. Let  $\kappa$  be as in (42),  $S_0 = \{j : \beta_j^* \neq 0\}$ ,  $\lambda_0 > 0$ ,  $0 < \eta < 1$ ,  $0 < \gamma_0 < 1/\kappa$ ,  $A > 1$ , and  $\xi \geq (A+1)/(A-1)$ . Suppose*

$$\lambda_0 \{1 + A/(1 - \kappa\gamma_0)\} \leq F_0(\xi, S; \phi_0) \eta e^{-\eta}, \quad F_* \leq F_2(\xi, S), \quad (45)$$

*for all  $S \supseteq S_0$  with  $|S \setminus S_0| \leq \ell^*$ , where  $F_0(\xi, S; \phi_0)$  is as in (22) and  $F_2(\xi, S)$  as in (43). Let  $\tilde{\beta}$  be an initial estimator of  $\beta$  and  $\hat{\beta}$  be as in (2) with  $\hat{w}_j = \dot{\rho}_\lambda(|\tilde{\beta}_j|)/\lambda$  and  $\lambda = A\lambda_0/(1 - \kappa\gamma_0)$ . Then,*

$$|\hat{\beta} - \beta^*|_2 \leq \frac{e^\eta}{F_*} \left\{ |\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2 + \left( \kappa + \frac{1}{\gamma_0 A} - \frac{\kappa}{A} \right) |\tilde{\beta} - \beta^*|_2 \right\}$$

*in the event  $\{|\tilde{\beta} - \beta\}_{S_0^c}|_2^2 \leq \gamma_0^2 \lambda^2 \ell^*\} \cap \{|z - \dot{\psi}(\beta^*)|_\infty \leq \lambda_0\}$ . Moreover, if (44) holds and  $\lambda_0 = \sigma \sqrt{(2/n) \log(2p/\epsilon_0)}$  with  $0 < \epsilon_0 < 1$ , then  $\mathbb{P}_{\beta^*} \{|z - \dot{\psi}(\beta^*)| \geq \lambda_0\} \leq \epsilon_0$ .*

Theorem 3 raises the possibility that  $\hat{\beta}$  improves  $\tilde{\beta}$  under proper conditions. Thus it is desirable to repeatedly apply this adaptive Lasso in the following way,

$$\hat{\beta}^{(k+1)} = \arg \min_{\beta} \left\{ \ell(\beta) + \sum_{j=1}^p \dot{\rho}_\lambda(\hat{\beta}_j^{(k)}) |\beta_j| \right\}, \quad k = 0, 1, \dots \quad (46)$$

Such multistage algorithms have been considered in [FL01, ZL08, Zha10b]. As discussed in Remark 4.1 below, it is beneficial to use a concave penalty  $\rho_\lambda$  in (46). Natural choices of  $\rho_\lambda$  include the smoothly clipped absolute deviation and minimax concave penalties [FL01, Zha10a].

**Theorem 4** *Let  $\{\kappa, S_0, \lambda_0, \eta, \gamma_0, A, \xi, \ell^*, \lambda\}$  be the same as Theorem 3. Let  $\widehat{\beta}^{(0)}$  be the unweighted Lasso with  $\widehat{w}_j = 1$  in (2) and  $\widehat{\beta}^{(\ell)}$  be the  $\ell$ -th iteration of the recursion (46) initialized with  $\widehat{\beta}^{(0)}$ . Let  $F_0(\xi, S_0; \phi_2)$  be the simple GIF in (22) with  $\phi_2(h) = |h|_2/|S|^{1/2}$ . Suppose (45) holds and*

$$e^\eta \{1 + (1 - \kappa\gamma_0)/A\} \sqrt{|S_0|}/F_0(\xi, S_0; \phi_2) \leq \gamma_0 \sqrt{\ell^*}. \quad (47)$$

Define  $r_0 = (e^\eta/F_*)\{\kappa + 1/(\gamma_0 A) - \kappa/A\}$ . Suppose  $r_0 < 1$ . Then,

$$|\widehat{\beta}^{(\ell)} - \beta^*|_2 \leq \frac{|\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2}{e^{-\eta}F_*(1 - r_0)/(1 - r_0^\ell)} + \frac{r_0^\ell e^\eta \lambda \{1 + (1 - \kappa\gamma_0)/A\}}{F_0(\xi, S_0; \phi_2)/|S_0|^{1/2}} \quad (48)$$

in the event

$$\left\{ |z - \dot{\psi}(\beta^*)|_\infty \leq \lambda_0 \right\} \cap \left\{ \frac{|\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2}{e^{-\eta}F_*(1 - r_0)} \leq \gamma_0 \lambda \sqrt{\ell^*} \right\}. \quad (49)$$

Moreover, if (44) holds and  $\lambda_0 = \sigma \sqrt{(2/n) \log(4p/\epsilon_0)}$  with  $0 < \epsilon_0 < 1$ , then the intersection of the events (49) and  $\{|\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2 \leq n^{-1/2} \sigma \sqrt{2|S_0| \log(4|S_0|/\epsilon_0)}\}$  happens with at least  $P_{\beta^*}$  probability  $1 - \epsilon_0$ , provided that

$$\frac{|\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + n^{-1/2} \sigma \sqrt{2|S_0| \log(4|S_0|/\epsilon_0)}}{e^{-\eta}F_*(1 - r_0)} \leq \frac{\gamma_0 A \lambda_0 \sqrt{\ell^*}}{1 - \kappa\gamma_0}. \quad (50)$$

**Remark 4.1** Define  $R^{(0)} = e^\eta \lambda \{1 + (1 - \kappa\gamma_0)/A\} |S_0|^{1/2}/F_0(\xi, S_0; \phi_2)$  and

$$R^{(\infty)} = \frac{|\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2}{e^{-\eta}F_*(1 - r_0)}, \quad R^{(\ell)} = (1 - r_0^\ell) R^{(\infty)} + r_0^\ell R^{(0)},$$

as in the right-hand side of (48). Theorem 4 asserts that  $|\widehat{\beta}^{(\ell)} - \beta^*| \leq 2R^{(\infty)}$  after  $\ell = |\log r_0|^{-1} \log(R^{(\infty)}/R^{(0)})$  iterations of the recursion (46). Under condition (44),

$$E_{\beta^*} R^{(\infty)} \leq \{|\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + 2\sigma \sqrt{|S_0|/n}\} e^\eta / \{F_*(1 - r_0)\}.$$

Suppose  $\rho_\lambda(t)$  is concave in  $t$ , then  $|\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 \leq \dot{\rho}_\lambda(0+) |S_0|^{1/2} = \lambda |S_0|^{1/2}$ . This component of  $E_{\beta^*} R^{(\infty)}$  matches the noise inflation due to model selection since  $\lambda \asymp \lambda_0 = \sigma \sqrt{(2/n) \log(p/\epsilon_0)}$ . This noise inflation diminishes when  $\min_{j \in S_0} |\beta_j^*| \geq \gamma \lambda$  when  $\dot{\rho}_\lambda(t) = 0$  for  $|t| \geq \gamma \lambda$ , yielding the super-efficient error bound  $E_{\beta^*} R^{(\infty)} \leq \{2\sigma \sqrt{|S_0|/n}\} e^\eta / \{F_*(1 - r_0)\}$ . This risk bound  $R^{(\infty)}$  is comparable with those for concave penalized least squares in linear regression [Zha10a].

**Remark 4.2** For  $\log(p/n) \asymp \log p$ , the penalty level  $\lambda$  in Theorems 3 and 4 are comparable with the best proven results and of the smallest possible order in linear regression. For  $\log(p/n) \ll \log p$ , the proper penalty level is expected to be of the order  $\sigma \sqrt{(2/n) \log(p/|S_0|)}$  under a vectorized sub-Gaussian condition which is slightly stronger than (44). This refinement for smaller  $p$  is beyond the scope of this paper.

**Remark 4.3** The constant factors used in Theorems 3 and 4 provide conditions of slightly weaker form than those based on sparse eigenvalues, although they typically do not imply each other due to differences in the dimension of covered models and various constant factors. If  $\phi_0(b) = M_3 |b_S|_1$  can be used as in (35), then  $M_3 |S| F_0(\xi, S; \phi_0) \geq F_2(\xi, S)$ . In GLM,  $\phi_0 = M_2 |b|_2$  can be used as in (32) to weaken this regularity condition. Since  $|b_S|_1 \leq |S|^{1/2} |b_S|_2$  and  $S_0 \subset S$ ,  $F_0(\xi, S_0; \phi_2) \geq F_0(\xi, S; \phi_2) \geq F_2(\xi, S)$ .

**Remark 4.4** Although Theorem 3 is valid for the smaller  $\xi \geq (A+1 - \kappa\gamma_0)/(A-1)$ , the proof of Theorem 4 requires  $\xi \geq (A+1)/(A-1)$ .

## 5 Selection consistency

In this section, we provide a selection consistency theorem for the  $\ell_1$  penalized convex minimization estimator, including both the weighted and unweighted cases. Let  $\|M\|_\infty = \max_{|u|_\infty \leq 1} |Mu|_\infty$  for matrices  $M$ .

**Theorem 5** Let  $\hat{\beta}$  be as in (2),  $\beta^*$  be a target vector,  $z_k^*$  be as in (5),  $\Omega_0$  in (6),  $S = \{j : \beta_j^* \neq 0\}$  and  $F(\xi, S; \phi_0, \phi)$  as in (17).

(i) Let  $0 < \eta \leq \eta^* \leq 1$  and  $\mathcal{B}_0^* = \{\beta : \phi_0(\beta - \beta^*) \leq \eta\}$ . Suppose

$$\sup_{\beta \in \mathcal{B}_0^*} \|W_{S^c}^{-1} \ddot{\psi}_{S^c, S}(\beta) \{\ddot{\psi}_S(\beta)\}^{-1} W_S\|_\infty \leq \kappa_0 < 1, \quad (51)$$

$$\sup_{\beta \in \mathcal{B}_0^*} \|W_{S^c}^{-1} \ddot{\psi}_{S^c, S}(\beta) \{\ddot{\psi}_S(\beta)\}^{-1}\|_\infty \leq \kappa_1. \quad (52)$$

Then,  $\{j : \hat{\beta}_j \neq 0\} \subseteq S$  in the event

$$\Omega_1^* = \Omega_0 \cap \left\{ |w_S|_\infty \lambda + z_0^* \leq \eta e^{-\eta} F(0, S; \phi_0, \phi_0), \kappa_1 z_0^* + z_1^* \leq (1 - \kappa_0) \lambda \right\}. \quad (53)$$

(ii) Let  $0 < \eta \leq \eta^* \leq 1$  and  $\mathcal{B}_0 = \{\beta : \phi_0(\beta - \beta^*) \leq \eta, \text{sgn}(\beta) = \text{sgn}(\beta^*)\}$ . Suppose (51) and (52) hold with  $\mathcal{B}_0^*$  replaced by  $\mathcal{B}_0$  and

$$\sup_{\beta \in \mathcal{B}_0} \|\{\ddot{\psi}_S(\beta)\}^{-1}\|_\infty \leq M_0, \quad (54)$$

Then,  $\text{sgn}(\hat{\beta}) = \text{sgn}(\beta^*)$  in the event

$$\Omega_1^* \cap \left\{ |w_S|_\infty \lambda + z_0^* < M_0^{-1} \min_{j \in S} |\beta_j^*| \right\}. \quad (55)$$

(iii) Suppose conditions of Theorem 2 hold for the GLM. Then, the conclusions of (i) and (ii) hold under the respective conditions if  $F(0, S; \phi_0, \phi_0)$  is replaced by  $F^*(\xi, S; M_2 |\cdot|_2)$  or  $F_*(\xi, S)$  or  $\kappa_*^2(\xi, S)/(M_3 |S|)$  with the respective  $\phi_0$  in Theorem 2.

For  $w_j = 1$ , this result is somewhat more specific in the radius  $\eta$  for the uniform irrerepresentable condition (51), compared with a similar extension of the selection consistency theory to the graphical Lasso by [RWR08]. In linear regression (10),  $\ddot{\psi}(\beta) = \Sigma = X'X/n$  does not depend on  $\beta$ , so that Theorem 5 with the special  $w_j = 1$  matches the existing selection consistency theory for the unweighted Lasso [MB06, Tro06, ZY06, Wai09]. We discuss below the  $\ell_1$  penalized logistic regression as a specific example.

**Example 5.1 (Logistic regression: selection consistency)** Suppose  $w_j = 1 = |x_j|_2^2/n$  where  $x_j$  are the columns of  $X$ . If (53) and (55) hold with  $z_0^*$  and  $z_1^*$  replaced by  $\sqrt{(\log(p/\epsilon_0))/(2n)}$ , then the respective conclusions of Theorem 5 hold with at least probability  $1 - \epsilon_0$  in  $P_{\beta^*}$ .

## 6 The sparsity of the Lasso and SRC

The results in Sections 2 and 3 are concerned with the estimation and prediction properties of  $\widehat{\beta}$ , but not dimension reduction. In this section, we provide upper bound for the dimension of  $\widehat{\beta}$ . For this purpose, we need to strengthen (21) to

$$e^{-\phi_0(b)}\Sigma^* \leq \ddot{\psi}(\beta^* + b) \leq e^{\phi_0(b)}\Sigma^*, \quad \forall b \in \mathcal{C}(\xi, S), \quad \phi_0(b) \leq \eta^*. \quad (56)$$

We assume the following sparse Riesz condition, or SRC [ZH08]:

$$c_* \leq u' \ddot{\psi}_A(\beta^*) u \leq c^*, \quad \frac{|S|}{2(1-\alpha)} \left( \frac{e^{2\eta} c^*}{c_*} + 1 - \alpha \right) \leq d^* \quad (57)$$

for certain constants  $\{c_*, c^*\}$ , integer  $d^*$ ,  $0 < \alpha \leq 1$ ,  $0 < \eta \leq \eta^* \leq 1$ , all  $A \supset S$  with  $|A| = d^*$  and all  $u \in \mathbb{R}^A$  with  $|u| = 1$ . The following theorem is an extension of the dimension bounds in [Zha10a] from linear regression.

**Theorem 6** *Let  $\beta^*$  and  $S$  be as in Theorem 1. Consider the  $\widehat{\beta}$  defined in (2) with  $w_j = 1$  for all  $j$ . Suppose (56) and (57) hold. Then,*

$$\#\{j : \widehat{\beta}_j \neq 0, j \notin S\} \leq d_1 = \left\lfloor \frac{|S|}{2(1-\alpha)} \left( \frac{e^{2\eta} c^*}{c_*} - 1 \right) \right\rfloor$$

in the event  $\Omega_1$  is defined in (18), provided that

$$\max_{A \supset S, |A| \leq d_1} |(\Sigma^*)_A^{-1/2} \dot{\ell}_A(\beta^*)|_2 \leq e^{-\eta} \alpha \lambda \sqrt{(d_1 - |S|)/c^*}.$$

For GLM, the results on the dimension bounds of the Lasso can be slightly simplified. Let  $\lambda_\xi = (\xi - 1)\lambda/(\xi + 1)$ . Suppose (56) and (57) hold and  $(\lambda + \lambda_\xi) \leq M_1 \eta e^{-\eta} F_*(0, S)$  with  $0 < \eta \leq 1$ . Then,

$$\#\{j : \widehat{\beta}_j \neq 0, j \notin S\} \leq d_1 = \left\lfloor \frac{|S|}{2(1-\alpha)} \left( \frac{e^{2\eta} c^*}{c_*} - 1 \right) \right\rfloor$$

in the event  $\{z^* < \lambda_\xi\}$ . The probability of the event  $\{z^* < \lambda_\xi\}$  can be calculated using Lemma 2 as in the previous sections.

## 7 Discussion

In this paper, we studied the estimation, prediction, selection and sparsity properties of the weighted  $\ell_1$ -penalized estimators in a general convex loss formulation.

We applied our general results to several important statistical models, including linear regression and generalized linear models. For linear regression, we extend the existing results to weighted/adaptive Lasso. For the GLMs, the  $\ell_q, q \geq 1$  error bounds for a general  $q \geq 1$  for the GLMs are not available in the literature, although  $\ell_1$  and  $\ell_2$  bounds have been obtained under different sets of conditions respectively in [vdG08, NRWY10]. Our fixed-sample analysis provides explicit constant factors in an explicit neighborhood of a target. Our oracle inequalities yields even sharper results for multistage recursive application of an adaptive Lasso.

An interesting aspect of the approach taken in this paper in dealing with general convex losses such as those for the GLM is that the conditions imposed on the Hessian naturally ‘converge’ to those for the linear regression as the convex loss ‘converges’ to a quadratic form.

A key quantity used in the derivation of the results is the generalized invertibility factor (17), which grow out of the idea of the  $\ell_2$  restricted eigenvalue but improves upon it. The use of GIF yields sharper bounds on the estimation and prediction errors. This was discussed in detail in the context of linear regression in [vdGB09, YZ10].

We assume that the convex function  $\psi(\cdot)$  is twice differentiable. Although this assumption is satisfied in many important and widely used statistical models, it would be interesting to extend the results obtained in this paper to models with less smooth loss functions, such as those in quantile regression and support vector machine.

## 8 Appendix

**Proof of Lemma 1.** Since  $\dot{\psi}(\hat{\beta}) - \dot{\psi}(\beta^*) = z - \dot{\psi}(\beta^*) - g$ , (3) implies

$$\Delta(\hat{\beta}, \beta^*) = \langle \hat{\beta}, z - \dot{\psi}(\beta^*) \rangle - \lambda |\widehat{W}\hat{\beta}|_1 - \langle \beta^*, z - \dot{\psi}(\beta^*) - g \rangle$$

and  $|g_j| \leq \hat{w}_j \lambda$ . Thus, (7) follows from  $|(z - \dot{\psi}(\beta^*))_j| \leq \hat{w}_j \lambda$  and  $\hat{w}_j \leq w_j$  in  $S$  in  $\Omega_0$ .

For (8), we have  $h_{S^c} = \hat{\beta}_{S^c}$  and  $\beta_{S^c}^* = 0$ , so that in  $\Omega_0$  (3) gives

$$\Delta(\hat{\beta}, \beta^*) = \langle \hat{\beta}_{S^c}, \{z - \dot{\psi}(\beta^*)\}_{S^c} \rangle - \lambda |\widehat{W}_{S^c} \hat{\beta}_{S^c}|_1 - \langle h_S, \{z - \dot{\psi}(\beta^*) - g\}_S \rangle$$

$$\begin{aligned}
&\leq |W_{S^c} \widehat{\beta}_{S^c}|_1(z_1^* - \lambda) + \langle h_S, g_S - \{z - \dot{\psi}(\beta^*)\}_S \rangle \\
&\leq |W_{S^c} \widehat{\beta}_{S^c}|_1(z_1^* - \lambda) + |h_S|_1(z_0^* + |w_S|_\infty \lambda).
\end{aligned}$$

This gives (8). Since  $\Delta(\widehat{\beta}, \beta^*) > 0$ ,  $h \in \mathcal{C}(\xi, S)$  when  $(|w_S|_\infty \lambda + z_0^*)/(\lambda - z_1^*) \leq \xi$ . For  $j \notin S$ ,  $h_j(\dot{\psi}(\beta + h) - \dot{\psi}(\beta))_j = \widehat{\beta}_j(z - \dot{\psi}(\beta^*) - g)_j \leq |\widehat{\beta}_j|(w_j \lambda - g_j) \leq 0$ .  $\square$

**Proof of Theorem 1.** Let  $h = \widehat{\beta} - \beta^*$ . Since  $\psi(\beta)$  is a convex function,

$$t^{-1} \Delta(\beta^* + th, \beta^*) = \frac{\partial}{\partial t} \left\{ \psi(\beta^* + th) - t \langle h, \dot{\psi}(\beta^*) \rangle \right\}$$

is an increasing function of  $t$ . For  $0 \leq t \leq 1$  and in the event  $\Omega_1$ , (8) implies

$$t^{-1} \Delta(\beta^* + th, \beta^*) \leq \Delta(h + \beta^*, \beta^*) < (|w_S|_\infty \lambda + z_0^*) |h_S|_1.$$

By (9) and (17),  $F(\xi, S; \phi_0, \phi_0) \leq \Delta(\beta^* + th, \beta^*) e^{\phi_0(th)} / \{t |h_S|_1 \phi_0(th)\}$  for  $\phi_0(th) \leq \eta^*$ . Thus, for  $\phi_0(th) \leq \min\{\eta^*, \phi_0(h)\}$  and in the event  $\Omega_1$ ,

$$\phi_0(th) e^{-\phi_0(th)} \leq \frac{\Delta(\beta^* + th, \beta^*)}{t |h_S|_1 F(\xi, S; \phi_0, \phi_0)} < \frac{|w_S|_\infty \lambda + z_0^*}{F(\xi, S; \phi_0, \phi_0)} \leq \eta e^{-\eta}.$$

If  $\eta^* < \phi_0(h)$ , the above inequality at  $\phi_0(th) = \eta^*$  would give  $\eta^* e^{-\eta^*} < \eta e^{-\eta}$ , which contradicts to  $\eta \leq \eta^* \leq 1$ . Thus,  $\eta^* \geq \phi_0(h)$  and  $\phi_0(th) e^{-\phi_0(th)} \leq \eta e^{-\eta}$  for all  $0 \leq t \leq 1$ . This implies  $\phi_0(h) \leq \eta \leq \eta^*$ . Another application of (8) yields

$$\phi(h) \leq \frac{\Delta(\beta^* + h, \beta^*) e^{\phi_0(h)}}{F(\xi, S; \phi_0, \phi) |h_S|_1} \leq \frac{(|w_S|_\infty \lambda + z_0^*) e^\eta}{F(\xi, S; \phi_0, \phi)}.$$

We obtain (20) by applying (19) with  $\phi = \phi_{1,S}$  to the right-hand side of (8).  $\square$

**Proof of Lemma 2.** (i) Since  $\dot{\psi}(\beta) = \sum_{i=1}^n x^i \dot{\psi}_0(x^i \beta) / n$  by (29),

$$\begin{aligned}
E_\beta \exp \left\{ \frac{n}{\sigma^2} b'(z - \dot{\psi}(\beta)) \right\} &= \exp \left[ \sum_{i=1}^n \frac{\psi_0(x^i(\beta + b)) - \psi_0(x^i \beta) - (x^i b) \dot{\psi}_0(x^i \beta)}{\sigma^2} \right] \\
&= \exp \left[ \sum_{i=1}^n \int_0^1 \frac{(x^i b)^2 \ddot{\psi}_0(x^i(\beta + tb))}{\sigma^2} (1-t) dt \right]. \quad (58)
\end{aligned}$$

This and (30) imply that for  $M_1 |Xb|_\infty \leq \eta_0$ ,

$$E_{\beta^*} \exp \left\{ \frac{n}{\sigma^2} b'(z - \dot{\psi}(\beta^*)) \right\} \leq \exp \left[ \frac{n e^{\eta_0} \langle b, \Sigma^* b \rangle}{2\sigma^2} \right]. \quad (59)$$



Since  $\max_{k=0,1} z_k^*/\lambda_k = \max_j t_j^{-1} |z_j - \dot{\psi}_j(\beta^*)|$  by (5),

$$\begin{aligned} \mathbb{P}_{\beta^*} \left\{ \max_{k=0,1} z_k^*/\lambda_k > 1 \right\} &\leq \sum_{j=1}^p \mathbb{P}_{\beta^*} \left\{ |z_j - \dot{\psi}_j(\beta^*)| > t_j \right\} \\ &\leq \sum_{j=1}^p \mathbb{E}_{\beta^*} \exp \left\{ \frac{n}{\sigma^2} b_j |z_j - \dot{\psi}_j(\beta^*)| - \frac{n}{\sigma^2} b_j t_j \right\} \end{aligned}$$

with  $b_j = e^{-\eta_0} t_j / \Sigma_{jj}^*$ . Since  $M_1 \max_{ij} |x_{ij}| b_j \leq \eta_0$ , (59) gives

$$\mathbb{P}_{\beta^*} \left\{ \max_{k=0,1} z_k^*/\lambda_k > 1 \right\} \leq \sum_{j=1}^p 2 \exp \left( - \frac{n e^{-\eta_0} t_j^2}{2 \sigma^2 \Sigma_{jj}^*} \right).$$

(ii) If (37) holds, we simply replace  $\ddot{\psi}_0(x^i(\beta + tb))$  by  $c_0$  in (58). The rest is simpler and omitted.  $\square$

**Proof of Theorem 2.** (i) Since  $F^*(\xi, S; \phi)$  in (32) is a lower bound of  $F(\xi, S; \phi_0, \phi)$  in (17), (39) follows from Theorem 1 with  $\phi_0(b) = M_2 |b|_2$ . The probability statement follows from Lemma 2. (ii) Since (21) holds for the  $\phi_0(b)$  in (33), we are allowed to use  $F_*(\xi, S) = F_0(\xi, S; \phi_0)$  in Corollary 1. The condition  $\eta^* = \infty$  is used since  $\phi_0(b)$  does not control  $M_1 |Xb|_\infty$ . (iii) We are also allowed to use  $\phi_0(b) = M_3 |b_S|_1$  in (35) due to  $M_1 |Xb|_\infty \leq \phi_0(b)$ .  $\square$

**Proof of Theorem 3.** Let  $h = \hat{\beta} - \beta^*$ ,  $w_j = \hat{w}_j$  and  $S = \{j : |\hat{\beta}_j| > \gamma_0 \lambda\} \cup S_0$ . For  $j \notin S$ ,  $w_j \lambda = \dot{\rho}_\lambda(\tilde{\beta}_j) \geq \dot{\rho}_\lambda(0+) - \kappa \gamma_0 \lambda = (1 - \kappa \gamma_0) \lambda$ , so that  $z_1^* = |W_{S^c}^{-1} \{z - \dot{\psi}(\beta^*)\}_{S^c}|_\infty \leq \lambda_0 / (1 - \kappa \gamma_0) = \lambda / A$ . We also have  $z_0^* \leq (1 - \kappa \gamma_0) \lambda / A$ . Since  $|\hat{w}|_\infty \leq 1$ , these bounds for  $z_0^*$  and  $z_1^*$  yield

$$\frac{|\hat{w}_S|_\infty \lambda + z_0^*}{\lambda - z_1^*} \leq \frac{\lambda + (1 - \kappa \gamma_0) \lambda / A}{\lambda - \lambda / A} = \frac{A + 1 - \kappa \gamma_0}{A - 1} \leq \xi.$$

Thus, by Lemma 1

$$h \in \mathcal{C}(\xi, S), \quad \Delta(\beta^* + h, \beta^*) \leq |h_S|_2 (|\hat{w}_S|_2 \lambda + |\{z - \dot{\psi}(\beta^*)\}_S|_2)$$

Since  $|S \setminus S_0| \leq |(\tilde{\beta} - \beta^*)_{S_0^c}|_2^2 / \gamma_0^2 \lambda^2 \leq \ell^*$ , we have

$$|w_S|_\infty \lambda + z_0^* \leq \lambda + (1 - \kappa \gamma_0) \lambda / A \leq F_0(\xi, S; \phi_0) \eta e^{-\eta}.$$

Thus,  $\phi_0(h) \leq \eta$  by (23). It follows that  $\Delta(\widehat{\beta}, \beta^*) \geq e^{-\eta} h' \Sigma h$  by (21), so that by (43),

$$e^{-\eta} F_* |h|_2 \leq e^{-\eta} F_2(\xi, S) |h|_2 \leq h' \Sigma h e^{-\eta} / |h_S|_2 \leq \Delta(\beta^* + h, \beta^*) / |h_S|_2$$

when  $|h_S| \neq 0$ . Consequently,

$$e^{-\eta} F_* |h|_2 \leq |\widehat{w}_S|_2 \lambda + |\{z - \dot{\psi}(\beta^*)\}_S|_2. \quad (60)$$

Since  $\widehat{w}_j \lambda = \dot{\rho}_\lambda(|\tilde{\beta}_j|) \leq \dot{\rho}_\lambda(|\beta_j^*|) + \kappa |\tilde{\beta}_j - \beta_j^*|$ , we have

$$|\widehat{w}_S|_2 \lambda \leq |\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + \kappa |\tilde{\beta} - \beta^*|_2.$$

Since  $|z - \dot{\psi}(\beta^*)|_\infty \leq (1 - \kappa \gamma_0) \lambda / A$ ,

$$\begin{aligned} |\{z - \dot{\psi}(\beta^*)\}_S|_2 &\leq |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2 + |S \setminus S_0|^{1/2} (1 - \kappa \gamma_0) \lambda / A \\ &\leq |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2 + |\tilde{\beta} - \beta^*|_2 (1 - \kappa \gamma_0) / (\gamma_0 A). \end{aligned}$$

Inserting the above inequalities into (60), we find that

$$e^{-\eta} F_* |\tilde{\beta} - \beta^*|_2 \leq |\dot{\rho}_\lambda(|\beta_{S_0}^*|)|_2 + |\{z - \dot{\psi}(\beta^*)\}_{S_0}|_2 + \left( \kappa + \frac{1}{\gamma_0 A} - \frac{\kappa}{A} \right) |\tilde{\beta} - \beta^*|_2.$$

The probability statement follows directly from (44) with the union bound.  $\square$

**Proof of Theorem 4.** Let  $R^{(\ell)}$  be as in Remark 4.1. For  $|z - \dot{\psi}(\beta^*)|_\infty \leq \lambda_0$ , Corollary 1 gives

$$|\widehat{\beta}^{(0)} - \beta^*|_2 \leq e^\eta (\lambda + \lambda_0) |S_0|^{1/2} / F_0(\xi, S_0; \phi_2) = R^{(0)}.$$

Under conditions (47) and (49), we have  $R^{(\ell)} \leq \gamma_0 \lambda \sqrt{\ell^*}$  for all  $\ell \geq 0$ . We prove (48) by induction. We have already proved (48) for  $\ell = 0$ . For  $\ell \geq 1$ , we let  $\tilde{\beta} = \widehat{\beta}^{(\ell-1)}$  and apply Theorem 3:  $|\widehat{\beta}^{(\ell)} - \beta^*|_2 \leq (1 - r_0) R^{(\infty)} + r_0 R^{(\ell-1)} = R^{(\ell)}$ . The probability statement follows directly from (44) with the union bound.  $\square$

**Proof of Theorem 5.** We first prove the more complicated part (ii). Let  $\tilde{z} = z - \dot{\psi}(\beta^*)$  and  $\lambda$  be fixed. Consider

$$\widehat{\beta}(\lambda, t) = \arg \min_{\beta} \left\{ \psi(\beta) - \langle \beta, \dot{\psi}(\beta^*) + t \tilde{z} \rangle + t \lambda \sum_{j=1}^p \widehat{w}_j |\beta_j| : \beta_{S^c} = 0 \right\} \quad (61)$$

as an artificial path for  $0 \leq t \leq 1$ . For each  $t$ , the KKT conditions for  $\widehat{\beta}(\lambda, t)$  are

$$g_j(\lambda, t) \begin{cases} = t\widehat{w}_j \lambda \operatorname{sgn}(\widehat{\beta}_j(\lambda, t)) & \forall \widehat{\beta}_j(\lambda, t) \neq 0 \\ \in t\widehat{w}_j[-1, 1], & \forall j \in S, \end{cases}$$

where  $g(\lambda, t) = -\dot{\psi}(\widehat{\beta}(\lambda, t)) + \dot{\psi}(\beta^*) + t\widetilde{z}$ . Let  $h(\lambda, t) = \widehat{\beta}(\lambda, t) - \beta^*$ . Since  $h_{S^c} = 0$ , the proof of Theorem 1 for  $\xi = 0$  yields

$$\phi_0(\widehat{\beta}(\lambda, t) - \beta^*) \leq \eta, \quad \forall 0 < t \leq 1. \quad (62)$$

Since  $\ddot{\psi}_S(\beta^*)$  is positive-definite,  $\widehat{\beta}(\lambda, 0+) = \beta^*$ . It follows that  $\operatorname{sgn}(\widehat{\beta}_S(\lambda, t)) = \operatorname{sgn}(\beta_S^*)$  for  $0 < t < t_1$  for a certain  $0 < t_1 \leq 1$ . An application of the differentiation operator  $D = (\partial/\partial t)$  to the KKT condition yields

$$\widetilde{z}_j - \ddot{\psi}_{j,S}(\widehat{\beta}(\lambda, t))\{(D\widehat{\beta})(\lambda, t)\}_S = \widehat{w}_j \lambda \operatorname{sgn}(\beta_j^*), \quad \forall j \in S, 0 < t < t_1.$$

Thus, for  $0 < t < t_1$

$$(D\widehat{\beta})_S(\lambda, t) = \{\ddot{\psi}_S(\widehat{\beta}(\lambda, t))\}^{-1}\{\widetilde{z}_S - \lambda\widehat{W}_S \operatorname{sgn}(\beta_S^*)\} \quad (63)$$

and with an application of the chain rule,

$$\begin{aligned} D\dot{\ell}_{S^c}(\widehat{\beta}(\lambda, t)) &= \ddot{\ell}_{S^c,S}(\widehat{\beta}(\lambda, t))(D\widehat{\beta})_S(\lambda, t) \\ &= \ddot{\psi}_{S^c,S}(\widehat{\beta}(\lambda, t))\{\ddot{\psi}_S(\widehat{\beta}(\lambda, t))\}^{-1}\{\widetilde{z}_S - \lambda\widehat{W}_S \operatorname{sgn}(\beta_S^*)\}. \end{aligned} \quad (64)$$

By (62),  $\widehat{\beta}(\lambda, t) \in \mathcal{B}_0$  for  $0 < t < t_1$ . It follows from (63), (54) and (55) that

$$|(D\widehat{\beta})_S(\lambda, t)|_\infty \leq M_0|\widetilde{z}_S - \lambda\widehat{W}_S \operatorname{sgn}(\beta_S^*)|_\infty \leq M_0(|\widehat{w}_S|_\infty \lambda + z_0^*) < \min_{j \in S} |\beta_j^*| - \epsilon_1$$

for  $0 < t < t_1$  and some  $\epsilon_1 > 0$ . Thus,  $|h_S(\lambda, t)|_\infty \leq tM_0(|w_S|_\infty \lambda + z_0^*) < \min_{j \in S} |\beta_j^*| - \epsilon_1$ . This implies  $\operatorname{sgn}(\widehat{\beta}(\lambda, t-)) = \operatorname{sgn}(\beta^*)$  for  $0 < t \leq 1$  by the continuity of  $\widehat{\beta}(\lambda, t)$  in  $t$ , i.e.  $t_1 = 1$ . Since  $|W_S^{-1}\widehat{W}_S v_S|_\infty \leq |v_S|_\infty$  for all  $v \in \mathbb{R}^p$  in  $\Omega_0$ , (64), (51) and (52) implies that for  $0 < t < 1$

$$\begin{aligned} |W_{S^c}^{-1}D\dot{\ell}_{S^c}(\widehat{\beta}(\lambda, t))| &\leq |W_{S^c}^{-1}\ddot{\psi}_{S^c,S}(\widehat{\beta}(\lambda, t))\{\ddot{\psi}_S(\widehat{\beta}(\lambda, t))\}^{-1}\widetilde{z}_S|_\infty \\ &\quad + \lambda|W_{S^c}^{-1}\ddot{\psi}_{S^c,S}(\widehat{\beta}(\lambda, t))\{\ddot{\psi}_S(\widehat{\beta}(\lambda, t))\}^{-1}W_S \operatorname{sgn}(\beta_S^*)|_\infty \\ &\leq \kappa_1 z_0^* + \kappa_0 \lambda. \end{aligned}$$

This implies  $|W_{S^c}^{-1}\dot{\ell}_{S^c}(\widehat{\beta}(\lambda, 1))|_\infty \leq \kappa_1 z_0^* + \kappa_0 \lambda + |W_{S^c}^{-1}\dot{\ell}_{S^c}(\beta^*)|_\infty \leq \kappa_1 z_0^* + z_1^* + \kappa_0 \lambda \leq \lambda$ . It follows that

$$\begin{cases} \dot{\ell}_j(\widehat{\beta}(\lambda, 1-)) = \widehat{w}_j \lambda \operatorname{sgn}(\beta_j^*), & \operatorname{sgn}(\beta_j^*) = \operatorname{sgn}(\widehat{\beta}(\lambda, 1-)), & j \in S \\ \dot{\ell}_j(\widehat{\beta}(\lambda, 1-)) \in \widehat{w}_j \lambda [-1, 1], & & j \notin S. \end{cases}$$

These are the KKT conditions for  $\widehat{\beta}(\lambda, 1-)$  with  $\operatorname{sgn}(\widehat{\beta}(\lambda, 1-)) = \operatorname{sgn}(\beta^*)$ .

The proof for part (ii) is similar, with  $\operatorname{sgn}(\beta^*)$  replaced by  $\operatorname{sgn}(\widehat{\beta}(\lambda, t))$  in the proof of part (i). Finally, in part (iii),  $F_0(\xi, S; \phi_0, \phi_0)$  is simply replaced by its lower bounds with the respective  $\phi_0$ .  $\square$

**Proof of Theorem 6.** Let  $A_1 = \{j : |g_j| = \lambda\} \cup S$ ,  $A_0 = A_1 \setminus S$  and  $\widehat{\Sigma} = \int_0^1 \ddot{\psi}(\beta^* + th)dt$ , where  $g$  is the negative gradient in (3) and  $h = \widehat{\beta} - \beta^*$ . Let  $g_{(A)} = (g_j I\{j \in A\})'$ . Consider the case  $|A_1| \leq d^*$  (e.g. with sufficiently small  $z^*$ ). Since  $g_{A_0} h_{A_0} = \lambda |h_{A_0}|_1$  by (3) and  $\{\dot{\ell}(\beta^* + h) - \dot{\ell}(\beta^*)\}_{A_1} = \widehat{\Sigma}_{A_1} h_{A_1}$ ,

$$g_{(A_0)} \widehat{\Sigma}_{A_1}^{-1} g_{(A_1)} = -g_{(A_0)} \widehat{\Sigma}_{A_1}^{-1} \dot{\ell}_{A_1}(\beta^* + h) \leq -\lambda |h_{A_0}|_1 + |\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_0)}| |\widehat{\Sigma}_{A_1}^{-1/2} \dot{\ell}_{A_1}(\beta^*)|.$$

Since  $|\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_0)}|_2^2 + |\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_1)}|_2^2 = |\widehat{\Sigma}_{A_1}^{-1/2} g_{(S)}|_2^2 + 2g_{(A_0)} \widehat{\Sigma}_{A_1}^{-1} g_{(A_1)}$ , we have

$$|\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_0)}|_2^2 + |\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_1)}|_2^2 \leq |\widehat{\Sigma}_{A_1}^{-1/2} g_{(S)}|_2^2 + 2|\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_0)}| |\widehat{\Sigma}_{A_1}^{-1/2} \dot{\ell}_{A_1}(\beta^*)|.$$

Thus, in the event  $|\widehat{\Sigma}_{A_1}^{-1/2} \dot{\ell}_{A_1}(\beta^*)| \leq \alpha \lambda \sqrt{|A_0|/(c^* e^\eta)}$  with  $0 < \alpha < 1$ , we have

$$(1 - \alpha) |\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_0)}|_2^2 + |\widehat{\Sigma}_{A_1}^{-1/2} g_{(A_1)}|_2^2 \leq |\widehat{\Sigma}_{A_1}^{-1/2} g_{(S)}|_2^2 + \alpha \lambda^2 |A_0|/(c_* e^\eta).$$

Since the eigenvalues of  $\widehat{\Sigma}_{A_1}$  lie in the interval  $c_* e^{-\eta}$  and  $c^* e^\eta$  and  $g_{A_0} = \lambda \operatorname{sgn}(\widehat{\beta}_{A_0})$ ,

$$\frac{(1 - \alpha) \lambda^2 |A_0|}{c^* e^\eta} + \frac{\lambda^2 |A_0| + |g_S|_2^2}{c_* e^\eta} \leq \frac{|g_S|_2^2}{c_* e^{-\eta}} + \frac{\alpha \lambda^2 |A_0|}{c^* e^\eta}.$$

This gives

$$2(1 - \alpha) |A_0| \leq \left( \frac{c^* e^{2\eta}}{c_*} - 1 \right) \frac{|g_S|_2^2}{\lambda^2} \leq \left( \frac{c^* e^{2\eta}}{c_*} - 1 \right) |S|.$$

We note that  $|\widehat{\Sigma}_{A_1}^{-1/2} \dot{\ell}_{A_1}(\beta^*)| \leq e^{\eta/2} \max_{A \supset S, |A| \leq d^*} |(\Sigma^*)_A^{-1/2} \dot{\ell}_A(\beta^*)|$ . We complete the proof by considering the artificial path (61).  $\square$

## References

- [Bre67] L. M. Bregman, *The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming*, USSR Computational Mathematics and Mathematical Physics **7** (1967), 200–217.
- [BRT09] Peter Bickel, Yaacov Ritov, and Alexandre Tsybakov, *Simultaneous analysis of Lasso and Dantzig selector*, Annals of Statistics **37** (2009), no. 4, 1705–1732.
- [BTW07] Florentina Bunea, Alexandre Tsybakov, and Marten H. Wegkamp, *Sparsity oracle inequalities for the Lasso*, Electronic Journal of Statistics **1** (2007), 169–194.
- [CDS98] S. Chen, D. L. Donoho, and M.A. Saunders, *Atomic decomposition by basis pursuit*, SIAM J. Sci. Comput. **20** (1998), 33–61.
- [CT05] Emmanuel J. Candes and Terence Tao, *Decoding by linear programming*, IEEE Trans. on Information Theory **51** (2005), 4203–4215.
- [CT07] E. Candes and T. Tao, *The dantzig selector: statistical estimation when  $p$  is much larger than  $n$  (with discussion)*, Annals of Statistics **35** (2007), 2313–2404.
- [FL01] Jianqing Fan and Runze Li, *Variable selection via nonconcave penalized likelihood and its oracle properties*, Journal of the American Statistical Association **96** (2001), 1348–1360.
- [GR04] E. Greenshtein and Y. Ritov, *Persistence in high-dimensional linear predictor selection and the virtue of overparametrization*, Bernoulli **10** (2004), 971–988.
- [HMZ08] J. Huang, S. Ma, and C.-H. Zhang, *Adaptive lasso for sparse high-dimensional regression models*, Statistica Sinica **18** (2008), 1603–1618.
- [Kol09] V. Koltchinskii, *The dantzig selector and sparsity oracle inequalities*, Bernoulli **15** (2009), 799–828.

- [MB06] Nicolai Meinshausen and Peter Bühlmann, *High-dimensional graphs and variable selection with the lasso*, Annals of Statistics **34** (2006), 1436–1462.
- [MB07] Lukas Meier and Peter Bühlmann, *Smoothing  $\ell_1$ -penalized estimators for high-dimensional time-course data*, Electronic Journal of Statistics **1** (2007), 597–615.
- [MN89] P. McCullagh and J.A. Nelder, *Generalized linear models*, Chapman & Hall, 1989.
- [MY09] N. Meinshausen and B. Yu, *Lasso-type recovery of sparse representations for high-dimensional data*, Annals of Statistics **37** (2009), 246–270.
- [NN07] Frank Nielsen and Richard Nock, *On the centroids of symmetrized bregman divergences*, CoRR **abs/0711.3242** (2007).
- [NRWY10] Sahand Negahban, Pradeep Ravikumar, Martin J. Wainwright, and Bin Yu, *A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizer*, Tech. Report arXiv:1010.2731, arXiv, 2010.
- [RBLZ08] Adam J. Rothman, Peter J. Bickel, Elizaveta Levina, and Ji Zhu, *Sparse permutation invariant covariance estimation*, Electronic Journal of Statistics **2** (2008), 494–515.
- [RWRY08] P. Ravikumar, M. J. Wainwright, G. Raskutti, and B. Yu, *Model selection in gaussian graphical models: High-dimensional consistency of  $\ell_1$ -regularized mle*, Advances in Neural Information Processing Systems (NIPS), vol. 21, 2008.
- [SBvdG10] N. Städler, P. Bühlmann, and S. van de Geer,  *$\ell_1$ -penalization for mixture regression models (with discussion)*, Test **19** (2010), no. 2, 209–285.
- [SZ11] Tungni Sun and Cun-Hui Zhang, *Scaled sparse linear regression*, Tech. Report arXiv:1104.4595, arXiv, 2011.

- [Tib96] R. Tibshirani, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society: Series B (Statistical Methodology) **58** (1996), 267–288.
- [Tro06] J. A. Tropp, *Just relax: convex programming methods for identifying sparse signals in noise*, IEEE Transactions on Information Theory **52** (2006), 1030–1051.
- [TT11] Ryan Tibshirani and Jonathan Taylor, *The solution path of the generalized lasso*, The Annals of Statistics **39** (2011), 1335–1371.
- [vdG07] S. van de Geer, *The deterministic lasso*, Tech. Report 140, ETH Zurich, Switzerland, 2007.
- [vdG08] ———, *High-dimensional generalized linear models and the lasso*, Annals of Statistics **36** (2008), 614–645.
- [vdGB09] S. van de Geer and P. Bühlmann, *On the conditions used to prove oracle results for the lasso*, Electronic Journal of Statistics **3** (2009), 1360–1392.
- [Wai09] M. J. Wainwright, *Sharp thresholds for noisy and high-dimensional recovery of sparsity using  $\ell_1$ -constrained quadratic programming (lasso)*, IEEE Transactions on Information Theory **55** (2009), 2183–2202.
- [YZ10] Fei Ye and Cun-Hui Zhang, *Rate minimaxity of the lasso and dantzig selector for the  $\ell_q$  loss in  $\ell_r$  balls*, Journal of Machine Learning Research **11** (2010), 3481–3502.
- [ZH08] Cun-Hui Zhang and Jian Huang, *The sparsity and bias of the Lasso selection in high-dimensional linear regression*, Annals of Statistics **36** (2008), no. 4, 1567–1594.
- [Zha09] Tong Zhang, *Some sharp performance bounds for least squares regression with  $L_1$  regularization*, Ann. Statist. **37** (2009), no. 5A, 2109–2144.
- [Zha10a] Cun-Hui Zhang, *Nearly unbiased variable selection under minimax concave penalty*, The Annals of Statistics **38** (2010), 894–942.

- [Zha10b] Tong Zhang, *Analysis of multi-stage convex relaxation for sparse regularization*, Journal of Machine Learning Research **11** (2010), 1087–1107.
- [Zha11] ———, *Adaptive forward-backward greedy algorithm for learning sparse representations*, IEEE Transactions on Information Theory (2011), to appear.
- [ZL08] Hui Zou and Runze Li, *One-step sparse estimates in nonconcave penalized likelihood models*, Annals of Statistics **36** (2008), no. 4, 1509–1533.
- [Zou06] Hui Zou, *The adaptive lasso and its oracle properties*, Journal of the American Statistical Association **101** (2006), 1418–1429.
- [ZY06] Peng Zhao and Bin Yu, *On model selection consistency of Lasso*, Journal of Machine Learning Research **7** (2006), 2541–2567.